# Deep Reinforcement Learning for Sequence-to-Sequence Models

Yaser Keneshloo, Tian Shi, Naren Ramakrishnan, Chandan K. Reddy, *Senior Member, IEEE*

arXiv:1805.09461v4 [cs.LG] 15 Apr 2019

*Abstract*—In recent times, sequence-to-sequence (seq2seq) models have gained a lot of popularity and provide state-of-the-art performance in a wide variety of tasks such as machine translation, headline generation, text summarization, speech to text conversion, and image caption generation. The underlying framework for all these models is usually a deep neural network comprising an encoder and a decoder. Although simple encoder-decoder models produce competitive results, many researchers have proposed additional improvements over these seq2seq models, e.g., using an attention-based model over the input, pointer-generation models, and self-attention models. However, such seq2seq models suffer from two common problems: 1) *exposure bias* and 2) *inconsistency between train/test measurement*. Recently, a completely novel point of view has emerged in addressing these two problems in seq2seq models, leveraging methods from reinforcement learning (RL). In this survey, we consider seq2seq problems from the RL point of view and provide a formulation combining the power of RL methods in decision-making with seq2seq models that enable remembering long-term memories. We present some of the most recent frameworks that combine concepts from RL and deep neural networks. Our work aims to provide insights into some of the problems that inherently arise with current approaches and how we can address them with better RL models. We also provide the source code for implementing most of the RL models discussed in this paper to support the complex task of abstractive text summarization and provide some targeted experiments for these RL models, both in terms of performance and training time.

*Index Terms*—Deep learning; reinforcement learning; sequence to sequence learning; Q-learning; actor-critic methods; policy gradients.

## I. INTRODUCTION

SEQUENCE-to-sequence (seq2seq) models constitute a common framework for solving sequential problems [1]. In seq2seq models, the input is a sequence of certain data units and the output is also a sequence of data units. Traditionally, these models are trained using a ground-truth sequence via a mechanism known as *teacher forcing* [2], where the teacher is the ground-truth sequence. However, due to some of the drawbacks of this training approach, there has been significant line of research connecting inference of these models with reinforcement learning (RL) techniques. In this paper, we aim to summarize such research in seq2seq training utilizing RL methods to enhance the performance of these models and discuss various challenges that arise when applying RL

Y. Keneshloo, T. Shi, N. Ramakrishnan, and C. K. Reddy are with the Discovery Analytics Center, Department of Computer Science at Virginia Tech, Arlington, VA. {yaserkl,tshi}@vt.edu, {naren,reddy}@cs.vt.edu. Corresponding author: yaserkl@vt.edu.
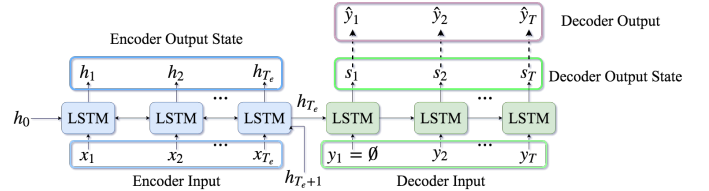
Fig. 1: A simple seq2seq model. The blue boxes correspond to the encoder part which has $T_e$ units. The green boxes correspond to the decoder part which has $T$ units.

methods to train a seq2seq model. We intend for this paper to provide a broad overview on the strength and complexity of combining seq2seq training with RL training and to guide researchers in choosing the right RL algorithm for solving their problem. In this section, we will briefly introduce the working of a simple seq2seq model and outline some of the problems that are inherent to seq2seq models. We will then provide an introduction to RL models and explain how these models could solve the problems of seq2seq models.

### A. Seq2seq Framework

Seq2seq models are common in various applications ranging from machine translation [3]–[8] , news headline generation [9], [10] , text summarization [11]–[14] , speech-to-text applications [15]–[18] , and image captioning [19]–[21].

In recent years, the general framework for solving these problems uses deep neural networks that comprise two main components: an encoder which reads the sequence of input data and a decoder which uses the output generated by the encoder to produce the sequence of final outputs. Fig 1 gives a schematic of this simple yet effective framework. The encoder and decoder are usually implemented by recurrent neural networks (RNN) such as Long Short-Term Memory (LSTM) [22]. The encoder takes a sequence of length $T_e$ inputs[1], $X = \{x_1, x_2, \cdots, x_{T_e}\}$, where $x_t \in \mathcal{A} = \{1, \cdots, |\mathcal{A}|\}$ is a single input coming from a range of possible inputs ($\mathcal{A}$), and generates the output state $h_t$. In addition, each encoder receives the the previous encoder's hidden state, $h_{t-1}$, and if the encoder is a bidirectional LSTM, it will also receive the state from the next encoder's hidden state, $h_{t+1}$, to generate its current hidden state $h_t$. The decoder, on the other hand, takes the last state from the encoder, i.e., $h_{T_e}$ and starts generating

---

[1]In this paper, we use input/output and action interchangeably since choosing the next input is akin to choosing the next action and generating the next output is akin to generating the next action.

an output of size $T < T_e$, $\hat{Y} = \{\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_T\}$, based on the current state of the decoder $s_t$ and the ground-truth output $y_t$. The decoder could also take as input an additional context vector $c_t$, which encodes the context to be used while generating the output [9]. The RNN learns a recursive function to compute $s_t$ and outputs the distribution over the next output:

$$
\begin{aligned}
h_{t'} &= \Phi_\theta(x_{t'}, h_t) \\
s_{t'} &= \Phi_\theta(y_t, s_t / h_{T_e}, c_t) \\
\hat{y}_{t'} &\sim \pi_\theta(y | \hat{y}_t, s_{t'})
\end{aligned}
\tag{1}
$$

where $t' = t + 1$, $\theta$ denotes the parameters of the model, and the function for $\pi_\theta$ and $\Phi_\theta$ depends on the type of RNN. A simple Elman RNN [23] would use a sigmoid function for $\Phi$ and a softmax function for $\pi$ [1]:

$$
\begin{aligned}
s_{t'} &= \sigma(W_1 y_t + W_2 s_t + W_3 c_t) \\
o_{t'} &= \text{softmax}(W_4 s_{t'} + W_5 c_t)
\end{aligned}
\tag{2}
$$

where $o_t$ is the output distribution of size $|\mathcal{A}|$ and the output $\hat{y}_t$ is selected from this distribution. $W_1$, $W_2$, $W_3$, $W_4$, and $W_5$ are matrices of learnable parameters of sizes $W_{1,2,3} \in R^{d \times d}$ and $W_{4,5} \in R^{d \times |\mathcal{A}|}$, where $d$ is the size of the input representation (e.g., size of the word embedding in text summarization). The input to the first decoder is a special input indicating the beginning of a sequence, denoted by $y_0 = \emptyset$ and the first forward hidden state $h_0$ and the last backward hidden state $h_{T_e+1}$ for the encoder are set to a zero vector. Moreover, the first hidden state for decoder $s_0$ is set to the output that is received from the last encoding state, i.e., $h_{T_e}$.

The most widely used method to train the decoder for sequence generation is called the teacher forcing algorithm [2], which minimizes the maximum-likelihood loss at each decoding step. Let us define $y = \{y_1, y_2, \cdots, y_T\}$ as the ground-truth output sequence for a given input sequence $X$. The maximum-likelihood training objective is the minimization of the following cross-entropy (CE) loss:

$$
\mathcal{L}_{CE} = -\sum_{t=1}^{T} \log \pi_\theta(y_t | y_{t-1}, s_t, c_{t-1}, X)
\tag{3}
$$

Once the model is trained with the above objective, the model generates an entire sequence as follows: Let $\hat{y}_t$ denotes the action (output) taken by the model at time $t$. Then, the next action is generated by:

$$
\hat{y}_{t'} = \arg\max_y \pi_\theta(y | \hat{y}_t, s_{t'})
\tag{4}
$$

This process could be improved by using beam search to find a reasonable good output sequence [7]. Now, given the ground-truth output $Y$ and the model generated output $\hat{Y}$, the performance of the model is evaluated with a specific measure. In seq2seq problems, discrete measures such as ROUGE [24], BLEU [25], METEOR [26], and CIDEr [27] are used to evaluate the model. For instance, ROUGE$_l$, an evaluation measure for textual seq2seq tasks, uses the largest common substring between $Y$ and $\hat{Y}$ to evaluate the goodness of the generated output. Algorithm 1 shows these steps.

---

**Algorithm 1** Training a simple seq2seq model

**Input**: Input sequences ($X$) and ground-truth output sequences ($Y$).
**Output**: Trained seq2seq model.
**Training Steps**:
**for** batch of input and output sequences $X$ and $Y$ **do**
    Run encoding on $X$ and get the last encoder state $h_{T_e}$.
    Run decoding by feeding $h_{T_e}$ to the first decoder and obtain the sampled output sequence $\hat{Y}$.
    Calculate the loss according to Eq. (3) and update the parameters of the model.
**end for**
**Testing Steps**:
**for** batch of input and output sequences $X$ and $Y$ **do**
    Use the trained model and Eq. (4) to sample the output $\hat{Y}$
    Evaluate the model using a performance measure, e.g., ROUGE
**end for**

---

### B. Problems with Seq2seq Models

One of the main issues with the current seq2seq models is that minimizing $\mathcal{L}_{CE}$ does not always produce the best results for the above discrete evaluation measures. Therefore, using cross-entropy loss for training a seq2seq model creates a mismatch in generating the next action during training and testing. As shown in Fig 1 and also according to Eq. (3), during training, the decoder uses the two inputs, the previous output state $s_{t-1}$ and the ground-truth input $y_t$, to calculate its current output state $s_t$ and uses it to generate the next action, i.e., $\hat{y}_t$. However, at the test time, as given in Eq. (4), the decoder completely relies on the previously generated action from the model distribution to predict the next action, since the ground-truth data is not available anymore. Therefore, in summary, the input to the decoder is from the ground-truth during training, but the input comes from the model distribution during model testing. This *exposure bias* [28] results in error accumulation during the output generation at test time, since the model has never been exclusively exposed to its own predictions during training. To avoid the *exposure bias* problem, we need to remove the ground-truth dependency during training and use only the model distribution to minimize Eq. (3). One way to handle this situation is through the scheduled sampling method [2] or Gibbs sampling [29]. In scheduled sampling, the model is first pre-trained using cross-entropy loss and will subsequently and slowly replace the ground-truth with a sampled action from the model. Therefore, a decision is randomly taken to whether use the ground-truth action with probability $\epsilon$, or an action coming from the model itself with probability $(1 - \epsilon)$. When $\epsilon = 1$, the model is trained using Eq. (3), and when $\epsilon = 0$ the model is trained based on the following loss:

$$
\mathcal{L}_{\text{Inference}} = -\sum_{t=1}^{T} \log \pi_\theta(\hat{y}_t | \hat{y}_1, \cdots, \hat{y}_{t-1}, s_t, c_{t-1}, X)
\tag{5}
$$

Note the difference between this equation and CE loss; in CE the ground-truth output $y_t$ is used to calculate the loss, while in Eq. (5), the output of the model $\hat{y}_t$ is used to calculate the loss.

Although scheduled sampling is a simple way to avoid the exposure bias, due to its random selection between choosing ground-truth output or model output, it does not provide a

clear solution for the back-propagation of error and therefore it is statistically inconsistent [30]. Recently, Goyal *et al.* [31] proposed a solution for this problem by creating a continuous relaxation over the argmax operation to create a differentiable approximation of the greedy search during the decoding steps.

As yet another line of research on avoiding the exposure bias problem, adversarial generative models are also proposed for various seq2seq models [32]–[35]. In general, adversarial models are comprised of a discriminator and a generator [36]. The generator tries to generate data similar to the ground-truth data while the discriminator's job is to discern whether the generated data is close to real data or it is a fake. Finally, the generator takes the feedback from the discriminator and optimizes its actions towards generating higher quality data. Since generator will only rely on its own output in generating the data, similar to the scheduled sampling, it is avoiding on reliance to the ground-truth data and hence avoids the exposure bias problem. However, adversarial generative models, in general, suffer from the reward sparsity [33], [35] and mode collapse [37] problems. Although, there are ways to avoid these two problems [38], studying these solutions is outside the scope of this work.

The second problem with seq2seq models is that, while the model training is done using the $\mathcal{L}_{CE}$, the model is typically evaluated during the test time using discrete and non-differentiable measures such as BLEU and ROUGE. This will create a mismatch between the training objective and the test objective and therefore could yield inconsistent results. Thus, a solution that could use these measures during training of the model will inherently solve this mismatch problem. Recently, it has been shown that both the *exposure bias* and non-differentiability of evaluation measures can be addressed by incorporating techniques from reinforcement learning [13], [39]–[41].

### C. Reinforcement Learning

In RL, a sequential Markov Decision Process (MDP) is considered, in which an agent interacts with an environment $\varepsilon$ over discrete time steps $t$ [42]. Let $M = (\mathcal{S}, \mathcal{A}, \mathcal{P}, R, s_0, \gamma, T)$ represent this discrete finite-horizon discounted Markov decision process, where $\mathcal{S}$ is the set of states, $\mathcal{A}$ is the set of actions, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}_+$ is the transition probability distribution, $\mathcal{R} : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$ is a reward function, $s_0 : \mathcal{S} \to \mathbb{R}_+$ is the initial state distribution, $\gamma \in [0, 1]$ a discount factor, and $T$ is the horizon.

The goal of the agent is to excel at a specific task, e.g., moving an object [43], [44], playing an Atari game [45], or generating news summary [13], [46]. The idea is that given the environment state at time $t$ as $s_t$, the agent picks an action $\hat{y}_t \in \mathcal{A}$, according to a (typically stochastic) policy $\pi(\hat{y}_t | s_t) : \mathcal{S} \times \mathcal{A} \to \mathbb{R}_+$ and observes a reward $r_t$ for that action[2]. The cumulative discounted sum of rewards is the objective function optimized by policy $\pi$. For instance, we can consider our seq2seq conditioned RNN as a stochastic policy that generates actions (selecting the next output) and

---

[2]we remove the subscript $t$ whenever it is clear from the context that we are in time $t$

receives the task reward based on a discrete measure like ROUGE as the return. The agent's goal is to maximize the expected discounted reward, $R_t = \mathbb{E}_\pi[\sum_{\tau=0}^{T} \gamma^\tau r_\tau]$, where the discounting factor $\gamma$ controls the trades off between the importance of immediate and future rewards. Under the policy $\pi$, we can define the values of the state-action pair $Q(s_t, y_t)$ and the state $V(s_t)$ as follows:

$$
\begin{aligned}
Q_\pi(s_t, y_t) &= \mathbb{E}[r_t | s = s_t, y = y_t] \\
V_\pi(s_t) &= \mathbb{E}_{y \sim \pi(s)}[Q_\pi(s_t, y = y_t)]
\end{aligned}
\tag{6}
$$

Note that the value function $V_\pi$ is defined over only the states whereas $Q_\pi$ is defined over (state, action) pairs. The $Q_\pi$ formulation is advantageous in model-free contexts since it can be applied to *current* states without having access to a model of the environment. In contrast, the $V_\pi$ formulation must, by necessity, be applied to *future* states and thus requires a model of the environment (i.e., which states and actions lead to which other future states). The preceding state-action function ($Q$ function for short) can be computed recursively with dynamic programming:

$$
Q_\pi(s_t, y_t) = \mathbb{E}_{s_{t'}}[r_t + \gamma \underbrace{\mathbb{E}_{y_{t'} \sim \pi(s_{t'})}[Q_\pi(s_{t'}, y_{t'})]}_{V_\pi(s_{t'})}]
\tag{7}
$$

Given the above definitions, we can define a function called advantage, denoted by $A_\pi$ relating the value function $V$ and $Q$ function as follows:

$$
\begin{aligned}
A_\pi(s_t, y_t) &= Q_\pi(s_t, y_t) - V_\pi(s_t) \\
&= r_t + \gamma \, \mathbb{E}_{s_{t'} \sim \pi(s_{t'} | s_t)}[V_\pi(s_{t'})] - V_\pi(s_t)
\end{aligned}
\tag{8}
$$

where $\mathbb{E}_{y \sim \pi(s)}[A_\pi(s, y)] = 0$ and for a deterministic policy, $y^* = \arg\max_y Q(s, y)$, it follows that $Q(s, y^*) = V(s)$, hence $A(s, y^*) = 0$. Intuitively, the value function $V$ measures how good the model could be when it is in a specific state $s$. The $Q$ function, however, measures the value of choosing a specific action when we are in such state. Given these two functions, we can obtain the advantage function which captures the importance of each action by subtracting the value of the state $V$ from the $Q$ function. In practice, seq2seq model is used as the policy which generates actions. The definition of an action, however, will be task-specific; e.g., for a text summarization task, the action denotes choosing the next token for the summary, whereas for a question answering task, the action might be defined as the start and end index of the answer in the document. Also, the definition of the reward function could vary from one application to another. For instance, in text summarization, measures like ROUGE and BLEU are commonly used while in image captioning, CIDEr and METEOR are common. Finally, the state of the model is usually defined as the decoder output state at each time step. Therefore, the decoder output state at each time is used as the current state of the model and is used to calculate our $Q$, $V$, and advantage functions. Table I summarizes the notations used in this paper.

### D. Paper Organization

In general, we define the following problem statement that we are trying to solve by combining these two different models of learning.

TABLE I: Notations used in this paper.

| Seq2seq Model Parameters | |
|---|---|
| $X$ | The sequence of input of length $T_e$, $X = \{x_1, x_2, \cdots, x_{T_e}\}$. |
| $Y$ | The sequence of ground-truth output of length $T$, $Y = \{y_1, y_2, \cdots, y_T\}$. |
| $\hat{Y}$ | The sequence of output generated by model of length $T$, $\hat{Y} = \{\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_T\}$. |
| $T_e$ | Length of the input sequence and number of encoders. |
| $T$ | Length of the output sequence and number of decoders. |
| $d$ | Size of the input and output sequence representation. |
| $\mathcal{A}$ | Input and output shared vocabulary. |
| $h_t$ | Encoder hidden state at time $t$. |
| $s_t$ | Decoder hidden state at time $t$. |
| $\pi_\theta$ | The seq2seq model with parameter $\theta$. |
| **Reinforcement Learning Parameters** | |
| $r_t = r(s_t, y_t)$ | The reward that the agent receives by taking action $y_t$ when the state of the environment is $s_t$ |
| $\hat{Y}$ | Sets of actions that the agent is taking for a period of time $T$, $\hat{Y} = \{\hat{y}_1, \hat{y}_2, \cdots, \hat{y}_T\}$. This is similar to the output that the seq2seq model is generating. |
| $\pi$ | The policy that the agent uses to take the action. |
| $\pi_\theta$ | Seq2seq models use RNNs with parameter $\theta$ for the policy. |
| $\gamma$ | Discount factor to reduce the effect of rewards from future actions. |
| $Q(s_t, y_t)$ $Q_\pi(s_t, y_t)$ | The $Q$-value (under policy $\pi$) that shows the estimated reward of taking action $y_t$ when at state $s_t$. |
| $Q_\Psi(s_t, y_t)$ | A function approximator with parameter $\Psi$ that estimates the $Q$-value given the state-action pair at time $t$. |
| $V(s_t)$ $V_\pi(s_t)$ | Value function which calculates the expectation of $Q$-value (under policy $\pi$) over all possible actions. |
| $V_\Psi(s_t)$ | A function approximator with parameter $\Psi$ that estimates the value function given the state at time $t$. |
| $A_\pi(s_t, y_t)$ | Advantage function (under policy $\pi$) which defines how good a state-action pair is w.r.t. the expected reward that can be received at this state. |
| $A_\Psi(s_t, y_t)$ | A function approximator with parameter $\Psi$ that estimates the advantage function for the state-action pair at time $t$. |

**Problem Statement**: *Given a series of input data and a series of ground-truth outputs, train a model that:*

- *Only relies on its own output, rather than the ground-truth, to generate the results (avoiding exposure bias).*
- *Directly optimize the model using the evaluation measure (avoiding mismatch between training and test measures).*

Although recently there had been a couple of survey articles on the topic of deep reinforcement learning [47], [48], these works heavily focused on the reinforcement learning methods and their applications in robotics and vision, while giving less emphasis to how these models could be used in a variety of other tasks. In this paper, we will summarize some of the most recent frameworks that attempted to find a solution for the above problem statement in a broad range of applications and explain how RL and seq2seq learning could benefit from each other in solving complex tasks. To this end, we will provide insights on some of the challenges and issues with the current models and how one can improve them with better RL models. The goal of this paper is to provide information about how we can broaden the power of seq2seq models with RL methods and understand challenges that exist in applying these methods to deep learning contexts. In addition, currently, there does not exist a good open-source framework for implementing these ideas. Along with this paper, we provide a library that combines state-of-the-art methods for the complex task of abstractive text summarization with recent techniques used in deep RL. The library provides a variety of different options and hyperparameters for training the seq2seq model using different RL models. Moreover, We provide experimental results on some of the most common techniques that are explained in this paper and we encourage researchers to experiment with other hyperparameters and explore how they can use this framework to gain better performance on

different seq2seq tasks. The contributions of this paper are summarized as follows:

- We provide a comprehensive summary of RL methods that are used in deep learning and specifically in the context of training seq2seq models.
- We summarize the challenges, advantages, and disadvantages of using different RL methods for seq2seq training.
- We provide guidelines on how one could improve a specific RL method to obtain a better and smoother training for seq2seq models.
- We provide an open-source library for implementing a complex seq2seq model using different RL techniques [3] along with experiments that aim for identifying an accurate estimate on the amount of improvement that RL algorithm provide for current seq2seq models.

This paper is organized as follows: Section III provides details of some of the common RL techniques used in training seq2seq models. We provide a brief introduction to different seq2seq models in Section IV and later explain various RL models that can be used along with the seq2seq model training process. We provide a summary of recent real-world applications that combine RL training with seq2seq training and in Section V we present the framework that we implemented and discuss the details about how this framework can applied to different seq2seq problems and provide experimental results for some of the well-known RL algorithm. Finally, in Section VI, we discuss conclusions of our work.

## II. SEQ2SEQ MODELS AND THEIR APPLICATIONS

Sequence to Sequence (seq2seq) models have been an integral part of many real-world problems. From Google Machine

[3] www.github.com/yaserkl/RLSeq2Seq/

Translation [4] to Apple's Siri speech to text [49], seq2seq models provide a clear framework to process information that is in the form of sequences. In a seq2seq model, the input and output are in the form of sequences of single units like sequence of words, images, or speech units. Table II provides a brief summary of various seq2seq models and their corresponding inputs and outputs. We also cite some of the important research papers for each application domain.

In recent years, different models and frameworks were proposed by researchers to achieve better and robust results on these tasks. For instance, attention-based models have been successfully applied to problems such as machine translation [3], text summarization [9], [10], question answering [64], image captioning [19], speech recognition [16], and object detection [89]. In an attention-based model, at each decoding step, the previous decoder output is combined with the information from the encoder's output at a specific position to select the best decoder output.

Although attention-based models can significantly improve the performance of seq2seq models in various tasks, in applications with large output space, it is challenging for the model to reach a desirable outcome.

On the other hand, there are more advanced models in seq2seq training like pointer-generator model [12], [90] and the transformers model which uses self-attention layers [91], but discussing these models is outside the scope of this paper.

Aside from these well-defined seq2seq problems, there are other related problems that partially work on the sequence of inputs but the output is not in the form of a sequence. Here are a few prominent applications that fall into this category.

- **Sentiment Analysis** [92]–[94]: The input is a sequence of words and the output is a single sentiment (positive, negative, or neutral).
- **Natural Language Inference** [95]–[97]: Given two sentences, one as a premise and the other as a hypothesis, the goal is to classify the relationship between these two sentences into one of the entailment, neutrality, and contradiction classes.
- **Sentiment Role Labeling** [98]–[101]: Given a sentence and a predicate, the goal is to answer questions like "who did what to whom and when and where".
- **Relation Extraction** [102]–[104]: Given a sentence, the goal is to identify whether a specific relationship exists in that sentence or not. For instance, based on the sentence "Barack Obama is married to Michelle Obama", we can extract the "spouse" relationship.
- **Pronoun Resolution** [105]–[108]: Given a sentence and a question about a pronoun in the sentence, the goal is to identify who that pronoun is referring to. For instance, in the sentence "Susan cleaned Alice's bedroom for the help **she** had given", the goal is to find who the word "she" is referring to.

Note that, although in these applications, only the input data is represented in terms of sequences, we still consider them to be seq2seq problems.

### A. Evaluation Measures

Seq2seq models are usually trained with cross-entropy loss, i.e., Eq. (3). However, the performance of these models is evaluated using discrete measures. There are various discrete measures that are used for evaluating these models and each application requires its own evaluation measure. We briefly provide a summary of these measures according to their application context:

- **ROUGE** [4] [24], **BLEU** [5] [25], **METEOR** [6] [26]: These are three of the most commonly used measures in applications such as machine translation, headline generation, text summarization, question answering, dialog generation, and other applications that require evaluation of text data. $ROUGE$ measure finds the common unigram ($ROUGE$-1), bigram ($ROUGE$-2), and largest common substring (LCS) ($ROUGE$-L) between the ground-truth text and the output generated by the model and calculate respective precision, recall, and F-score for each measure. $BLEU$ works similar to $ROUGE$ but through a modified precision calculation, it inclines to provide higher scores to outputs that are closer to human judgement. In a similar manner, $METEOR$ uses the harmonic mean of unigram precision and recall and it gives higher importance to recall than the precision. Although these methods are designed to work for all text-based applications, $METEOR$ is more often used in machine translation tasks, while $ROUGE$ and $BLEU$ are mostly used in text summarization, question answering, and dialog generation.
- **CIDEr** [7] [27], **SPICE** [8] [109]: $CIDEr$ is frequently used in image and video captioning tasks in which having captions that have higher human judgement scores is more important. Using sentence similarity, the notions of grammaticality, saliency, importance, accuracy, precision, and recall are inherently captured by these metrics.
  $SPICE$ is a recent evaluation metric proposed for image captioning that tries to solve some of the problems of $CIDEr$ and $METEOR$ by mapping the dependency parse trees of the caption to the semantic scene graph (contains objects, attributes of objects, and relations) extracted from the image. Finally, it uses the F-score that is calculated using the tuples of the generated and ground-truth scene graphs to provide the caption quality score.
- **Word Error Rate** (**WER**): This measure, which is mostly used in speech recognition, finds the number of substitutions, deletions, insertions, and corrections required to change the generated output to the ground-truth and combines them to calculate the $WER$.

### B. Datasets

In this section, we briefly describe some of the datasets that are commonly used in various seq2seq models. We provide a

---

[4] https://github.com/andersjo/pyrouge/
[5] https://www.nltk.org/_modules/nltk/translate/bleu_score.html
[6] http://www.cs.cmu.edu/~alavie/METEOR/
[7] https://github.com/vrama91/cider
[8] http://www.panderson.me/spice/

TABLE II: A summary of different applications of seq2seq models. In seq2seq models, the input and output are sequences of unit data. The input column provides information about the sequences of data that are fed into the model and the output column provides information about the sequences of data that the model generates as its output.

| Application | Problem Description | Input | Output | References |
|---|---|---|---|---|
| **Machine Translation** | Translating a sentence from a source language to a target language | A sentence (sequence of words) in language X (e.g., English) | Another sentence (sequence of words) in language Y (e.g., French) | [1], [6], [7] [3], [50], [51] |
| **Text Summarization Headline Generation** | Summarizing a document into a more concise and shorter text | A long document like a news article (sequence of words) | A short summary/headline (sequence of words) | [9], [10], [14], [52] [12], [53]–[58] |
| **Question Generation** | Generating interesting questions from a text document or an image | A piece of text (sequence of words) or image (sequence of layers) | A set of questions (sequence of words) related to the text or image | [59]–[62] |
| **Question Answering** | Given a text document or an image and a question, find the answer to the question | A textual question (sequence of words) or an image (sequence of layers) | A single word answer from a document or the start and end index of the answer in the document | [63]–[65] |
| **Dialogue Generation** | Generate a dialogue between two agents e.g., between a robot and human | A dialogue from the first agent (sequence of words) or audibles (sequence of speech units) | A dialogue from the second agent (sequence of words) or audibles (sequence of speech units) | [66]–[69] |
| **Semantic Parsing** | Generating automatic SQL queries from a given human-written description | A human-written description of the query (sequence of words) | The SQL command equivalent to that description | [70], [71] |
| **Image Captioning** | Given an image, generate a caption that explains the content of the image | An image (sequence of layers) | The caption (sequence of words) describing that image | [72]–[75] [19], [20], [76] |
| **Video Captioning** | Given a video clip, generate a caption that explains the content of the video | A video (sequence of images) | The caption (sequence of words) describing the video | [77]–[80] |
| **Computer Vision** | Finding interesting events in a video clip, e.g., predicting the next action of a specific object in the video | A video (sequence of images) | Differs from application to application. For instance, one might be interested in determining the next action of a specific object or entity in the video | [81]–[84] |
| **Speech Recognition** | Given a segment of audible input (e.g., speech), convert it to text and vice versa | A speech (sequence of speech units) | The text of the input speech (sequence of words) | [16], [17], [85], [86] |
| **Speech Synthesis** | Given a segment of text it generates its audible sounds | A text (sequence of words) | A speech representing (sequence of speech units) representing its audible sounds | [87], [88] |

short list of some of the most common datasets that are used in various seq2seq applications as follows:

- **Machine Translation**: The most common dataset used for Machine Translation task is the **WMT'14** [9] dataset which contains 850M words from English-French parallel corpora of UN (421M words), Europarl (61M words), news commentary (5.5M words), and two crawled corpora of 90M and 272.5M words. The data pre-processing for this dataset is usually done following the code [10] provided by Axelrod *et al.* [110] .

- **Text Summarization**: One of the main datasets used in text summarization is the **CNN-Daily Mail** dataset [111] which is part of the DeepMind Q&A Dataset [11] and contains around 287K news articles along with 2 to 4 highlights (summary) for each news article [12]. Recently, another dataset, called **Newsroom**, was released by Connected

Experiences Lab [13] [112] which contains 1.3M news articles and various metadata information such as the title and summary of the news. The document summarization challenge [14] also provides some datasets for text summarization. More specifically, in this dataset, **DUC-2003** and **DUC-2004** which contain 500 news articles (each paired with 4 different human-generated reference summaries) from the New York Times and Associated Press Wire services, respectively. Due to the small size of this dataset, researchers usually use this dataset only for evaluation purposes.

- **Headline Generation**: Headline generation is similar to the task of text summarization and typically all the datasets that are used in text summarization will be useful in headline generation, too. There is a big dataset which is called **Gigaword** [113] and contains more than 8M news articles from multiple news agencies like New York Times and Associate Press. However, this dataset is not freely available and researchers are required to buy the license to be able to use it though one can still find pre-trained models on

[9] http://www.statmt.org/wmt14/translation-task.html

[10] http://www-lium.univ-lemans.fr/~schwenk/cslm_joint_paper/

[11] https://cs.nyu.edu/~kcho/DMQA/

[12] For downloading and pre-processing please refer to: https://github.com/abisee/cnn-dailymail

[13] https://summari.es/

[14] https://duc.nist.gov/data.html

different tasks using this dataset [15].

- **Question Answering and Question Generation**: The **CNN-Daily Mail** dataset was originally designed for question answering and is one of the earliest datasets that is available for tackling this problem. However, recently two large-scale datasets that are solely designed for this problem were released. Stanford Question Answering Dataset (**SQuAD**) [16](1.0 and 2.0) [114], [115] is a dataset for reading comprehension and contains more than 100K pairs of questions and answers collected by crowd-sourcing over a set of Wikipedia articles. The answer to each question is a segment which identifies the start and end indices of the answer within the article. The second dataset is called **TriviaQA** [17] [116], and similar to **SQuAD**, it is designed for reading comprehension and question answering task. This dataset contains 650K triples of questions, answers, and evidences (which help to find the answer).

- **Dialogue Generation**: The dataset for this problem usually comprises of dialogues between different people. The **Open-Subtitles** dataset [18] [117], **Movie Dialog dataset** [19] [118], and **Cornell Movie Dialogues** Corpus [20] [119] are three examples of these types of datasets. **OpenSubtitles** contains conversations between movie characters for more than 20K movies in 20 languages. The **Cornell Movie Dialogues** corpus contains more than 220K dialogues between more than 10K movie characters.

- **Semantic Parsing**: Recently Zhong *et al.* [70] released a dataset called **WikiSQL** [21] for this problem which contains 80654 hand-annotated questions and SQL queries distributed across 24241 tables from Wikipedia. Although this is not the only dataset for this problem but it offers a larger set of examples from other datasets such as **WikiTableQuestion** [22] [120] and **Overnight** [121].

- **Sentiment Analysis**: For this application, **Amazon product review** [23] [122] dataset is one of the largest dataset which contains more than 82 million product reviews from May 1996 to July 2014 in its de-duplicated version. Another big dataset for this task is the **Stanford Sentiment Treebank (SSTb)** [24] [94], which includes fine grained sentiment labels for 215,154 phrases in the parse trees of 11,855 sentences.

- **Natural Language Inference**: **Stanford Natural Language Inference** (SNLI) [25] [123] is the standard dataset for this task which contains 570K human-written English sentence pairs manually labeled for the three classes entailment, contradiction, and neutral. The **Multi-Genre Natural Language Inference** (MultiNLI) [26] [124] corpus is another new dataset which is collected through crowd-sourcing and contains 433K sentence pairs annotated with textual entailment information.

- **Semantic Role Labeling**: **Proposition Bank** (PropBank) [27] [125] is the standard dataset for this task which contains a corpus of text annotation with information about basic semantic propositions in seven different languages.

- **Relation Extraction**: **Freebase** [28] [126] is a huge dataset containing billions of triples: the entity pair and the specific relationship between them which are selected from the New York Times corpus (NYT).

- **Pronoun Resolution**: The **OntoNotes 5.0** dataset [29] is the standard dataset for this task. Specifically, researchers use the Chinese portion of this dataset to do the pronoun resolution in Chinese [105], [106], [108].

- **Image Captioning**: There are two datasets that are mainly used in image captioning. The first one is the **COCO** dataset [30] [127] which is designed for object detection, segmentation, and image captioning. This dataset contains around 330K images amongst which 82K images are used for training and 40K used for validation in image captioning. Each image has five ground-truth captions. **SBU** [128] is another dataset which consists of 1M images from Flickr and contains descriptions provided by image owners when they uploaded the images to Flickr.

- **Video Captioning**: For this problem, **MSR-VTT** [31] [129] and **YouTube2Text/MSVD** [32] [130] are two of the widely used datasets. MSR-VTT consists 10K videos from a commercial video search engine each containing 20 human annotated captions and YouTube2Text/MSVD which has 1970 videos each containing on an average 40 human annotated captions.

- **Image Classification**: The most popular dataset in computer vision is the **MNIST** dataset [33] [131]. This dataset consists of handwritten digits and contains a training set of 60K examples and a test set of 10K examples. Aside from this dataset, there is a huge list of datasets that are used for various computer vision problems and explaining each of them is beyond the scope of this paper [34].

- **Speech Recognition**: **LibriSpeech ASR Corpus** [35] [132] is one of the main datasets used for the speech recognition task. This dataset is free and is composed of 1000 hours of segmented and aligned 16kHz English speech which is derived from audiobooks. **Wall Street Journal** (WSJ) also has two Continuous Speech Recognition corpora containing 70 hours of speech and text from a corpus of Wall Street Journal news text. However, unlike the LibriSpeech dataset, this dataset is not freely available and researchers have to buy a license to use it. Similar to the WSJ dataset, **TIMIT** [36]

---

[15] http://opennmt.net/Models/

[16] https://rajpurkar.github.io/SQuAD-explorer/

[17] http://nlp.cs.washington.edu/triviaqa/

[18] http://opus.nlpl.eu/OpenSubtitles.php

[19] http://fb.ai/babi

[20] http://www.cs.cornell.edu/~cristian/Cornell_Movie-Dialogs_Corpus.html

[21] https://github.com/salesforce/WikiSQL

[22] https://nlp.stanford.edu/software/sempre/wikitable/

[23] http://jmcauley.ucsd.edu/data/amazon/

[24] https://nlp.stanford.edu/sentiment/

[25] https://nlp.stanford.edu/projects/snli/

[26] https://www.nyu.edu/projects/bowman/multinli/

[27] http://propbank.github.io/

[28] https://old.datahub.io/dataset/freebase

[29] https://catalog.ldc.upenn.edu/LDC2013T19

[30] http://cocodataset.org/

[31] http://ms-multimedia-challenge.com/2017/challenge

[32] http://www.cs.utexas.edu/users/ml/clamp/videoDescription/

[33] http://yann.lecun.com/exdb/mnist/

[34] Please refer to this link for a comprehensive list of datasets that are used in computer vision: http://riemenschneider.hayko.at/vision/dataset/

[35] http://www.openslr.org/12/

[36] https://catalog.ldc.upenn.edu/ldc93s1

is another dataset containing the read speech data. It contains time-aligned orthographic, phonetic, and word transcriptions of recordings for 630 speakers of eight major dialects of American English in which each of them are reading ten phonetically sentences.

## III. Reinforcement Learning Methods

In reinforcement learning, the goal of an agent interacting with an environment is to maximize the expectation of the reward that it receives from the actions. Therefore, the focus is on maximizing one of the following objectives:

$$\text{maximize } \mathbb{E}_{\hat{y}_1,\cdots,\hat{y}_T \sim \pi_\theta(\hat{y}_1,\cdots,\hat{y}_T)}[r(\hat{y}_1,\cdots,\hat{y}_T)] \quad (9)$$

$$\underset{y}{\text{maximize }} A_\pi(s_t, y_t) \quad (10)$$

$$\underset{y}{\text{maximize }} A_\pi(s_t, y_t) \rightarrow Maximize_y \ Q_\pi(s_t, y_t) \quad (11)$$

There are various ways in which one can solve this problem. In this section, we explain the solutions in detail and provide their strengths and weaknesses. Different methods aim solving this problem by trying one of the following approaches: (i) solve this problem through Eq. (9); (ii) solve the expected discounted reward $\mathbb{E}[R_t = \sum_{\tau=t}^{T} \gamma^{\tau-t} r_\tau]$; (iii) solve it by maximizing the advantage function (Eq. (10)); and (iv) solve it by maximizing $Q$ function using Eq. (11). Most of these methods are suitable choices for improving the performance of seq2seq models, but depending on the approach that is chosen for training the reinforced model, the training procedure for seq2seq model also changes. The first and one of the simplest algorithms that will be discussed in this section is the Policy Gradient (PG) method which aims to solve Eq. (9). Section III-B discusses Actor-Critic (AC) methods which improve the performance of PG models by solving Eq. (10) through Eq. (7) expansion on $Q$-function. Section III-C discusses $Q$-learning models that aim at maximizing the $Q$ function (Eq. (11)) to improve the PG and AC models. Finally, Section III-D will provide more details about some of the recent models which improve the performance of $Q$-learning models.

### A. Policy Gradient

In all reinforcement algorithms, an agent takes some action according to a specific policy $\pi$. The definition of a policy varies according to the specific application that is being considered. For instance, in text summarization, the policy is a language model $p(y|X)$ that, given input $X$, tries to generate the output $y$. Now, let us assume that our agent is represented by an RNN and takes actions from a policy $\pi_\theta$[37]. In a deterministic environment, where the agent takes discrete actions, the output layer of the RNN is usually a softmax function and it generates the output from this layer. In Teacher Forcing, a set of ground-truth sequences are given and the actions are chosen according to the current policy during training and the reward is observed only at the end of the sequence or when an End-Of-Sequence (EOS) signal is seen. Once the agent reaches the end of sequence, it compares

[37]In seq2seq model, this represents $\pi_\theta(y_t|\hat{y}_{t-1}, s_t, c_{t-1})$ in Eq. (1)

the sequence of actions from the current policy $(\hat{y}_t)$ against the ground-truth action sequence $(y_t)$ and calculate a reward based on any specific evaluation metric. The goal of training is to find the parameters of the agent in order to maximize the expected reward. This loss is defined as the negative expected reward of the full sequence:

$$\mathcal{L}_\theta = -\mathbb{E}_{\hat{y}_1,\cdots,\hat{y}_T \sim \pi_\theta(\hat{y}_1,\cdots,\hat{y}_T)}[r(\hat{y}_1,\cdots,\hat{y}_T)] \quad (12)$$

where $\hat{y}_t$ is the action chosen by the model at time $t$ and $r(\hat{y}_1,\cdots,\hat{y}_T)$ is the reward associated with the actions $\hat{y}_1,\cdots,\hat{y}_T$. Usually, in practice, one will approximate this expectation with a single sample from the distribution of actions acquired by the RNN model. Hence, the derivative for the above loss function is given as follows:

$$\nabla_\theta \mathcal{L}_\theta = -\underset{\hat{y}_{1\ldots T} \sim \pi_\theta}{\mathbb{E}} [\nabla_\theta \log \pi_\theta(\hat{y}_{1\ldots T}) r(\hat{y}_{1\ldots T})] \quad (13)$$

Using the chain rule, this equation can be re-written as follows [133]:

$$\nabla_\theta \mathcal{L}_\theta = \frac{\partial \mathcal{L}_\theta}{\partial \theta} = \sum_t \frac{\partial \mathcal{L}_\theta}{\partial o_t} \frac{\partial o_t}{\partial \theta} \quad (14)$$

where $o_t$ is the input to the softmax function. The gradient of the loss $\mathcal{L}_\theta$ with respect to $o_t$ is given by [41], [133]:

$$\frac{\partial \mathcal{L}_\theta}{\partial o_t} = \Big(\pi_\theta(y_t|y_{t-1}, s_t, c_{t-1}) - \mathbf{1}(\hat{y}_t)\Big)(r(\hat{y}_1,\cdots,\hat{y}_T) - r_b) \quad (15)$$

where $\mathbf{1}(\hat{y}_t)$ is the 1-of-$|\mathcal{A}|$ representation of the ground-truth output and $r_b$ is a baseline reward and could be any value as long as it is not dependent on the parameters of the RNN network. This equation is very similar to the gradient of a multi-class logistic regression. In logistic regression, the cross-entropy gradient is the difference between the prediction and the actual 1-of-$|\mathcal{A}|$ representation of the ground-truth output:

$$\frac{\partial \mathcal{L}_\theta^{CE}}{\partial o_t} = \pi_\theta(y_t|y_{t-1}, s_t, c_{t-1}) - \mathbf{1}(y_t) \quad (16)$$

Note that, in Eq. (15), the generated output from the model is used as a surrogate ground-truth for the output distribution, while, in Eq. (16), the ground-truth is used to calculate the gradient.

The goal of the baseline reward is to force the model to select actions that yield a reward $r > r_b$ and discourage those that have reward $r < r_b$. Since only one sample is being used to calculate the gradient of the loss function, it is shown that, having this baseline would reduce the variance of the gradient estimator [41]. If the baseline is not dependent on the parameters of the model $\theta$, Eq. (15) is an unbiased estimator. To prove this, we simply need to show that adding the baseline reward $r_b$ does not have any effect on the expectation of loss:

$$\mathbb{E}_{\hat{y}_{1\ldots T} \sim \pi_\theta}[\nabla_\theta \log \pi_\theta(\hat{y}_{1\ldots T}) r_b] = r_b \sum_{\hat{y}_{1\ldots T}} \nabla_\theta \pi_\theta(\hat{y}_{1\ldots T})$$

$$= r_b \nabla_\theta \sum_{\hat{y}_{1\ldots T}} \pi_\theta(\hat{y}_{1\ldots T}) = r_b \nabla_\theta 1 = 0 \quad (17)$$

This algorithm is called $REINFORCE$ [41] and is a simple yet elegant policy gradient algorithm for seq2seq problems. One of the challenges with this method is that the model suffers from high variance since only one sample is
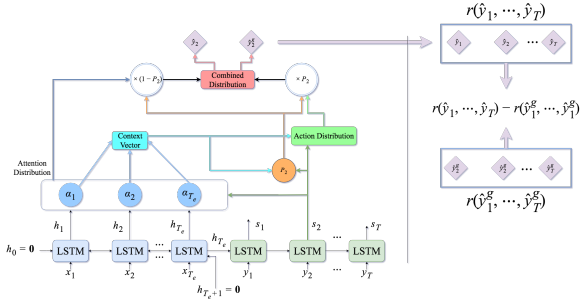
Fig. 2: A simple attention-based pointer-generation seq2seq model with Self-Critic reward. At each decoding step, the context vector for that decoder is calculated and combined with the decoder output to get the action distribution. In pointer-generation model, the attention distribution is further combined with the action distribution through switches called pointers to get the final distribution over the actions. From each output distribution, a specific action $\hat{y}_2$ is sampled and the greedy action $\hat{y}_2^g$ is extracted. The difference of the rewards from sampling and greedy sequence is used to update the loss function.

used for training at each time step. To alleviate this problem, at each training step, one can sample $N$ sequences of actions and update the gradient by averaging over all these $N$ sequences as follows:

$$\mathcal{L}_\theta = \frac{1}{N} \sum_{i=1}^{N} \sum_t \log \pi_\theta(\hat{y}_{i,t}|\hat{y}_{i,t-1}, s_{i,t}, c_{i,t-1}) \times \left( r(\hat{y}_{i,1}, \cdots, \hat{y}_{i,T}) - r_b \right) \tag{18}$$

Having this, the baseline reward could be set to be the mean of the $N$ rewards that are sampled, i.e., $r_b = 1/N \sum_{i=1}^{N} r(\hat{y}_{i,1}, \cdots, \hat{y}_{i,T})$. Algorithm 2 shows how this method works.

As another solution to reduce the variance of the model, **Self-Critic** (SC) models are proposed [40]. In these SC models, rather than estimating the baseline using current samples, the output of the model obtained by a greedy-search (the output at the time of inference) is used as the baseline. Hence, the sampled output of the model is used as $\hat{y}_t$ and greedy selection of the final output distribution is used for $\hat{y}_t^g$ where the superscript $g$ indicates greedy selection. Following this mechanism, the new objective for the REINFORCE model would become as follows:

$$\mathcal{L}_\theta = \frac{1}{N} \sum_{i=1}^{N} \sum_t \log \pi_\theta(\hat{y}_{i,t}|\hat{y}_{i,t-1}, s_{i,t}, c_{i,t-1}) \times \left( r(\hat{y}_{i,1}, \cdots, \hat{y}_{i,T}) - r(\hat{y}_{i,1}^g, \cdots, \hat{y}_{i,T}^g) \right) \tag{19}$$

Fig 2 shows how an attention-based pointer-generator seq2seq model can be used to extract the reward and its baseline in Self-Critic model.

The second problem with this method is that the reward is only observed after the full sequence of actions is sampled. This might not be a pleasing feature for most of the seq2seq models. If we see the partial reward of a given action at time $t$, and the reward is bad, the model needs to select a better action for the future to maximize the reward. However, in the REINFORCE algorithm, the model is forced to wait till the

---

**Algorithm 2** REINFORCE algorithm

**Input**: Input sequences ($X$), ground-truth output sequences ($Y$), and (preferably) a pre-trained policy ($\pi_\theta$).
**Output**: Trained policy with REINFORCE.
**Training Steps**:
**while** not converged **do**
    Select a batch of size $N$ from $X$ and $Y$.
    Sample $N$ full sequence of actions:
    $\{\hat{y}_1, \cdots, \hat{y}_T \sim \pi_\theta(\hat{y}_1, \cdots, \hat{y}_T)\}_1^N$.
    Observe the sequence reward and calculate the baseline $r_b$.
    Calculate the loss according to Eq. (18).
    Update the parameters of network $\theta \leftarrow \theta + \alpha \nabla_\theta \mathcal{L}_\theta$.
**end while**
**Testing Steps**:
**for** batch of input and output sequences $X$ and $Y$ **do**
    Use the trained model and Eq. (4) to sample the output $\hat{Y}$.
    Evaluate the model using a performance metric, e.g. $ROUGE_l$.
**end for**

---

end of the sequence to observe its performance. Therefore, the model often generates poor results or takes longer to converge. This problem magnifies especially in the beginning of the training phase where the model is initialized randomly and thus selects arbitrary actions. To alleviate this problem to a certain extent, Ranzato *et al.* [28] suggested to pre-train the model for a few epochs using the cross-entropy loss and then slowly switch to the REINFORCE loss. Finally, as another way to solve the high variance problem of the REINFORCE algorithm, importance sampling [134], [135] can also be used. The basic underlying idea of using the importance sampling with REINFORCE algorithm is that rather than sampling sequences from the current model, one can sample them from an old model and use them to calculate the loss.

### B. Actor-Critic Model

As mentioned in Section III-A, adding a baseline reward is a necessary component of the PG algorithm in order to reduce the variance of the model. In PG, the average reward from multiple samples in the batch was used as the baseline reward for the model. In Actor-Critic (AC) model, the goal is to train an estimator for calculating the baseline reward. For computing this quantity, AC models try to maximize the advantage function through Eq. (7) extension. Therefore, these methods are also called Advantage Actor-Critic (A2C) models.

In these models, the goal is to solve this problem using the following objective:

$$A_\pi(s_t, y_t) = Q_\pi(s_t, y_t) - V_\pi(s_t) = r_t + \gamma \mathbb{E}_{s_{t'} \sim \pi(s_{t'}|s_t)}[V_\pi(s_{t'})] - V_\pi(s_t) \tag{20}$$

Similar to the PG algorithm, to avoid the expensive inner expectation calculation, we can only sample once and approximate advantage function as follows:

$$A_\pi(s_t, y_t) \approx r_t + \gamma V_\pi(s_{t'}) - V_\pi(s_t) \tag{21}$$

Now, in order to estimate $V_\pi(s)$, a function approximator can be used to approximate the value function. In AC, neural networks is typically used as the function approximator for the value function. Therefore, we fit a neural network $V_\pi(s; \Psi)$ with parameters $\Psi$ to approximate the value function. Now, if we consider $r_t + \gamma V_\pi(s_{t'})$ as the expectation of reward-to-go at time $t$, $V_\pi(s_t)$ could play as a surrogate for the

baseline reward. Similar to the PG, the variance of the model would be high since only one sample is used to train the model. Therefore, the variance can be reduced using multiple samples. In the AC model, the Actor (our policy, $\theta$) provides samples (policy states at time $t$ and $t+1$) for the Critic (neural network estimating value function, $V_\pi(s; \Psi)$) and the Critic returns the estimation to the Actor, and finally, the Actor uses these estimations to calculate the advantage approximation and update the loss according to the following equation:

$$\mathcal{L}_\theta = \frac{1}{N} \sum_{i=1}^{N} \sum_t \log \pi_\theta(\hat{y}_i, |\hat{y}_{i,t-1}, s_{i,t}, c_{i,t-1}) A_\Psi(s_{i,t}, y_{i,t}) \tag{22}$$

Therefore, in the AC models, the inference at each time $t$ would be as follows:

$$\arg\max_y \pi_\theta(\hat{y}_t|\hat{y}_{t-1}, s_t, c_{t-1}) A_\Psi(s_t, y_t) \tag{23}$$

Fig. 3 provides an illustration of how this model works at one of the decoding steps.

*1) Training Critic Model:* As mentioned in the previous section, the Critic is a function estimator which tries to estimate the expected reward-to-go for the model at time $t$. Therefore, training the Critic is basically a regression problem. Usually, in AC models, a neural network is used as the function approximator and the value function is trained using the mean square error:

$$\mathcal{L}(\Psi) = \frac{1}{2} \sum_i ||V_\Psi(s_i) - v_i||^2 \tag{24}$$

where $v_i = \sum_{t'=t}^{T} r(s_{i,t'}, y_{i,t'})$ is the true reward-to-go at time $t$. During training the Actor model, we collect $(s_i, v_i)$ pairs and pass them to the Critic model to train the estimator. This model is called *on-policy AC*, which refers to the fact that the samples are collected at the current time to train the Critic model. However, the samples that are passed to the Critic will be correlated to each other which causes poor generalization for the estimator. These methods could be turn to off-policy by collecting training samples into a memory buffer and select mini-batches from this memory buffer and train the Critic network. Off-policy AC provides a better training due to avoiding the correlation of samples that exists in the on-policy methods. Therefore, most the models that we discuss in this paper are primarily off-policy and use a memory buffer for training the Critic model.

Algorithm 3 shows the batch Actor-Critic algorithm since for training the Critic network, we use a batch of state-rewards pair. In the online AC algorithm, the Critic network is simply updated using just one sample and, as expected, online AC algorithm has a higher variance due to reliance on one sample for training the network. To alleviate this problem for online AC, we can use Synchronous Advantage AC learning or Asynchronous Advantage AC (A3C) learning [136]. In the synchronous approach, $N$ different threads are used to train the model and each thread performs online AC for one sample and at the end of the algorithm, the gradient of these $N$ threads is used to update the gradient of the Actor model. In the more widely used A3C algorithm, as soon as a thread calculates $\theta$, it will send the update to other threads and other threads
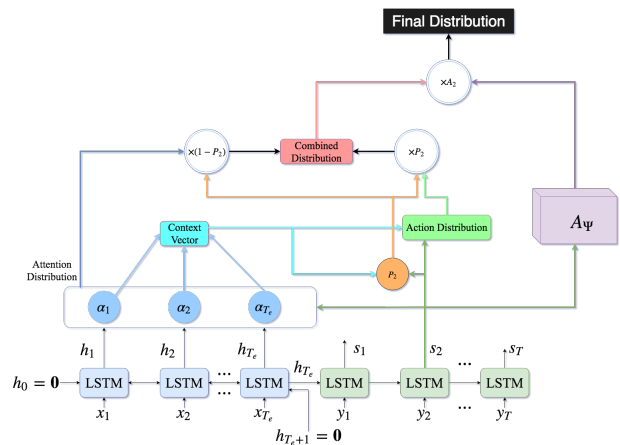


Fig. 3: A simple Actor-Critic model with an attention-based pointer-generation seq2seq model as the Actor. The Critic model is shown on the right side of the picture with a purple box. The purple box $A_\Psi$, which represents the Critic model, takes as input the decoder output at time $t = 2$, i.e., $s_2$, and estimate the advantage values through either (value function estimation, DQN, DDQN, or dueling net) for each action.

use the updated $\theta$ to train the model. A3C is an on-policy method with multi-step returns while there are other methods like Retrace [137], UNREAL [138], and Reactor [139] which provide the off-policy variations of this model by using the memory buffer. Also, ACER [140] mixes on-policy (from current run) and off-policy (from memory) to train the Critic network.

In general, AC models usually have low variance due to the batch training and the use of critic as the baseline reward, but they are not unbiased if the critic is erroneous and makes a lot of mistakes. As mentioned in Section III-A, PG algorithm has high variance but it provides an unbiased estimator. Now, if the PG and AC model are combined, we will likely be ending up with a model that has no bias and low variance. This idea comes from the fact that, for deterministic policies (like seq2seq models), a partially observable loss could be driven by using the $Q$-function as follows [39], [141]:

$$\begin{aligned} \mathcal{L}_\theta = \frac{1}{N} \sum_{i=1}^{N} \sum_t &\log \pi_\theta(\hat{y}_{i,t}|\hat{y}_{i,t-1}, s_{i,t}, c_{i,t-1}) \times \\ &\left( Q_\Psi(s_{i,t}) - V_{\Psi'}(s_{i,t}) \right) \end{aligned} \tag{25}$$

However, this model requires training two different networks for $Q_\Psi$ function and $V_{\Psi'}$ function as the baseline. Note that the same model cannot be used to estimate both $Q$ function and value function since the estimator will not be an unbiased estimator anymore [142]. As yet another solution to create a trade-off between the bias and variance in AC, Schulman *et al.* [143] proposed the Generalized Advantage Estimation ($GAE$) model as follows:

$$A_\Psi^{GAE}(s_t, y_t) = \sum_{i=t}^{T} (\gamma\lambda)^{i-t} \left( r(s_i, y_i) + \gamma V_\Psi(s_{i+1}) - V_\Psi(s_i) \right) \tag{26}$$

**Algorithm 3** Batch Actor-Critic Algorithm
***
**Input**: Input sequences ($X$), ground-truth output sequences ($Y$), and (preferably) a pre-trained Actor model ($\pi_\theta$).
**Output**: Trained Actor and Critic models.
**Training Steps**:
Initialize the Actor (Seq2seq) model, $\pi_\theta$.
Initialize the Critic (ValueNet) model, $V_\Psi$.
**while** not converged **do**
  **Training Actor**:
  Select a batch of size $N$ from $X$ and $Y$.
  Sample $N$ full sequences of actions based on the Actor.
  model, $\pi_\theta$.
  **for** $n = 1, \cdots, N$ **do**
    **for** $t = 1, \cdots, T$ **do**
      Calculate the true (discounted) reward-to-go:
      $v_t = \sum_{t'=t}^{T} \gamma^{t'-t} r(s_{i,t'}, y_{i,t'})$.
      Store training pairs for Critic: $(s_t, v_t)$.
    **end for**
  **end for**

  **Training Critic**:
  Select a batch of size $N_c$ from the pool of state-rewards pairs.
  collected from Actor.
  **for** $n = 1, \cdots, N_c$ **do**
    Collect the value estimates $\hat{v}_n$ from $V_\Psi$ for each
    state-rewards pair.
  **end for**
  Minimize the Critic loss using Eq. (24).

  **Updating Actor**:
  Use the estimated value for $V_\Psi(s_t)$ and $V_\Psi(s_{t'})$
  to calculate the loss using Eq. (22).
  Update parameters of the model using $\theta \leftarrow \theta + \alpha \nabla_\theta \mathcal{L}(\theta)$.
**end while**
***

where $\lambda$ controls the trade-off between the bias and variance such that large values of $\lambda$ yield to larger variance and lower bias, while small values of $\lambda$ do the opposite.

### C. Actor-Critic with Q-Learning

As mentioned in the previous section, the value function is used to maximize the advantage function. As an alternative to solve the maximization of advantage estimates, we can try to solve the following objective function:

$$Maximize_y \ A_\pi(s_t, y_t) \rightarrow Maximize_y \ Q_\pi(s_t, y_t) - \underbrace{V_\pi(s_t)}_{0} \quad (27)$$

This is true since we are trying to find the actions that maximize the advantage estimate and since value function does not rely on the actions, we can simply remove them from the maximization objective. Therefore, the advantage maximization problem is simplified to $Q$ function estimation problem. This method is called $Q$-learning and it is one of the most commonly used algorithms for RL problems. The $Q$-learning is called to a family of off-policy algorithm used to learn a $Q$-function. Similar to this method, SARSA algorithm [144] is an on-policy algorithm for calculating the $Q$-function. The major difference between SARSA and $Q$-Learning, is that the maximum reward for the next state is not necessarily used for updating the $Q$-values. In $Q$-learning, the Critic tries to provide an estimation for the $Q$-function. Therefore, given that the policy $\pi_\theta$ is being used, our goal is

to maximize the following loss at each training step:

$$\mathcal{L}_\theta = \frac{1}{N} \sum_{i=1}^{N} \sum_t \log \pi_\theta(\hat{y}_{i,t}|\hat{y}_{i,t-1}, s_{i,t}, c_{i,t-1}) Q_\Psi(s_{i,t}, y_{i,t}) \quad (28)$$

Similar to the value network training, the $Q$-function estimation is a regression problem and the Mean Squared Error (MSE) is used for training it. However, one of the differences between the $Q$-function training and value function training is the way in which the true estimates are chosen. In value function estimation, the ground-truth data is used to calculate the true reward-to-go as $v_i = \sum_{t'=t}^{T} r(s_{i,t'}, y_{i,t'})$, however, in $Q$-learning, the estimation from the network approximator itself is used to train the regression model:

$$\mathcal{L}(\Psi) = \frac{1}{2} \sum_i ||Q_\Psi(s_i, y_i) - q_i||^2$$
$$q_i = r_t + \gamma max_{y'} Q_\Psi(s'_i, y'_i) \quad (29)$$

where $s'_i$ and $y'_i$ are the state and action at the next time, respectively. Although the $Q$-value estimation has no direct relation to the true $Q$-values calculated using ground-truth data, in practice, it is known to provide good estimation and faster training due to not collecting ground-truth reward at each step of the training. However, there are no rigorous studies that analyze how far are these estimates from the true $Q$-values. As shown in Eq. (29), the true $Q$ estimations is calculated using the estimation from network approximator at time $t + 1$, i.e. $max_y' Q_\Psi(s'_i, y'_i)$. Although, not relying on the true ground-truth estimation and explicitly using the reward function might seems to be a bad idea, however in practice it is shown that these models provide better and more robust estimators. Therefore, the training process in $Q$-learning consists of first collecting a dataset of experiences $e_t = (s_t, y_t, s_{t'}, r_t)$ during training our Actor model and then use them to train the network approximator. This is the standard way of training the $Q$-network and was frequently used in earlier temporal-difference learning models. But, there is a problem with this method. Generally, the Actor-Critic models with neural network as function estimator are tricky to train and unless there are guarantees that the estimator is good, the model does not converge. Although the original $Q$-learning method is proven to converge [145], [146], when a neural networks is used to approximate the estimator, the convergence guarantee no longer holds. Usually, since samples are coming from a specific sets of sequences, there is a correlation between the samples that are chosen to train the model. Thus, this may cause any small updates to Q-network to significantly change the data distribution, and ultimately affects the correlations between $Q$ and the target values. Recently, Mnih *et al.* [45] proposed using an *experience buffer* [147][38] to store the experiences from different sequences and then randomly select a batch from this dataset and train the $Q$-network. Similar to the off-policy AC model, one benefit of using this buffer is the potential to increase efficiency of the model by re-using the experiences in multiple updates and reducing the variance of the model. Since, by sampling uniformly from the buffer, the correlation of samples used in the updates is reduced. As another improvement to the *experience buffer*, a prioritized version of this buffer is used in which, to select

***

[38]In some literatures, it is called a *replay buffer*

---

**Algorithm 4** Deep $Q$-Learning

---

**Input**: Input sequences ($X$), ground-truth output sequences ($Y$), and preferably a pre-trained Actor model ($\pi_\theta$).
**Output**: Trained Actor and Critic models.
**Training Steps**:
Initialize the Actor (Seq2seq) model, $\pi_\theta$.
Initialize the Critic ($Q$-Net) model, $Q_\Psi$.
**while** not converged **do**
    **Training Seq2seq Model**:
    Select a batch of size $N$ from $X$ and $Y$.
    Sample $N$ full sequences of actions based on the Actor model, $\pi_\theta$.
    **for** $n = 1, \cdots, N$ **do**
        **for** $t = 1, \cdots, T$ **do**
            Collect experience $e_t = (s_t, y_t, s_{t'}, r_t)$ and add them to the *experience buffer*.
        **end for**
    **end for**

    **Training $Q$-Net**:
    Select a batch of size $N_q$ from the *experience buffer*.
    based on the reward.
    **for** $n = 1, \cdots, N_q$ **do**
        Estimate $\hat{q_n} = Q_\Psi(s_n, y_n)$.
        Calculate the true estimation:
$$q_n = \begin{cases} r_n & s_n'\text{==EOS} \\ r_n + \gamma max_{y'} Q_\Psi(s_n', y_n') & \text{otherwise.} \end{cases}$$
        Store $(\hat{q}_n, q_n)$.
    **end for**
    **Updating $Q$-Net**:
    Minimize the loss using Eq. (29).
    Update the parameters of network, $\Psi$.

    **Updating Seq2seq Model**:
    Use the estimated $Q$ values for $\hat{q}_n = Q_\Psi(s_n, y_n)$.
    to calculate the loss using Eq. (28).
    Update parameters of the model using $\theta \leftarrow \theta + \alpha \nabla_\theta \mathcal{L}(\theta)$.
**end while**

---

**Algorithm 5** Double Deep $Q$-Learning

---

**Input**: Input sequences ($X$), ground-truth output sequences ($Y$), and preferably a pre-trained Actor model ($\pi_\theta$).
**Output**: Trained Actor and Critic models.
**Training Steps**:
Initialize the Actor (Seq2seq) model, $\pi_\theta$.
Initialize the two Critic models:
current $Q$-Net, $Q_\Psi$, and target $Q$-net, $Q_{\Psi'}$: $Q_{\Psi'} \leftarrow Q_\Psi$.
**while** not converged **do**
    **Training Seq2seq Model**:
    Select a batch of size $N$ from $X$ and $Y$.
    Sample $N$ full sequences of actions based on the Actor model, $\pi_\theta$.
    **for** $n = 1, \cdots, N$ **do**
        **for** $t = 1, \cdots, T$ **do**
            Collect experience $e_t = (s_t, y_t, s_{t'}, r_t)$ and add them to the *experience buffer*.
        **end for**
    **end for**

    **Training $Q$-Net**:
    Select a batch of size $N_q$ from the *experience buffer*
    based on the reward.
    **for** $n = 1, \cdots, N_q$ **do**
        Estimate $\hat{q_n} = Q_\Psi(s_n, y_n)$
        Calculate the true estimation:
$$q_n = \begin{cases} r_n & s_n'\text{==EOS} \\ r_n + \gamma Q_\Psi(s_n', \arg\max_{y_t'} Q_{\Psi'}(s_t', y_t')) & \text{otherwise.} \end{cases}$$
        Store $(\hat{q}_n, q_n)$.
    **end for**
    **Updating current $Q$-Net**:
    Minimize the loss using Eq. (29).
    Update the parameters of network, $\Psi$.

    **Updating target $Q$-Net every $N_u$ iterations**:
    $\Psi' \leftarrow \Psi$ or using Polyak averaging:
    $\Psi' \leftarrow \tau\Psi' + (1-\tau)\Psi, \tau = \frac{1000 - (\text{Current Step}\%1000)}{1000}$.

    **Updating Seq2seq Model**:
    Use the estimated $Q$-values for $\hat{q}_n = Q_\Psi(s_n, y_n)$
    to calculate the loss using Eq. (28).
    Update parameters of the model using $\theta \leftarrow \theta + \alpha \nabla_\theta \mathcal{L}(\theta)$.
**end while**

---

the mini-batches during training, only samples that have low temporal difference error are selected [148]. Algorithm 4 provides the pseudo-code for a $Q$-learning algorithm called Deep $Q$-Network or DQN.

### D. Advanced Q-Learning

*1) Double Q-Learning:* One of the problems with the Deep $Q$-Network (DQN) is the overestimation of $Q$-values as discussed in [149], [150]. Specifically, the problem lies in the fact that the ground-truth reward is not used to train these models and the same network is used to calculate both the estimation of network $Q_\Psi(s_i, y_i)$ and true values for regression training, $q_i$. To alleviate this problem, one could use two different networks in which the first one chooses the best action when calculating $max_{y'} Q_\Psi(s_n', y_n')$ and the other calculate the estimation of $Q$ value, i.e., $Q_\Psi(s_i, y_i)$. In practice, a modified version of the current DQN network is used as the second network in which the current network freezes its parameters for a certain period of time and updates the second network, periodically. Let us call the second network as the target network with parameters $\Psi'$. We know that $max_{y'} Q_\Psi(s_n', y_n')$ is the same as choosing the best action according to the network $Q_\Psi$. Therefore, this equation can be re-written as $Q_\Psi(s_t', \arg\max_{y_t'} Q_\Psi(s_t', y_t'))$. As shown in this equation, $Q_\Psi$ is used for both calculating the $Q$-value and finding the best action. Given a target network, the best action is chosen using our target network and the $Q$-value is

estimated using the current network. Hence, using the target network, $Q_{\Psi'}$, the $Q$-estimation will be given as follows:

$$q_t = \begin{cases} r_t & s_n'\text{==EOS} \\ r_t + \gamma Q_\Psi(s_t', \arg\max_{y_t'} Q_{\Psi'}(s_t', y_t')) & \text{otherwise.} \end{cases} \tag{30}$$

where EOS stands for the End-Of-Sequence action. This method is called Double DQN (DDQN) [149], [151] and is shown to resolve the problem of overestimation in DQN and provides more realistic estimations. However, even this model suffers from the fact that there is no relation between the true $Q$-values and the estimation provided by the network. Algorithm 5 shows the pseudocode for this model.

*2) Dueling Networks:* DDQN tried to solve one of the problems with DQN model by using two networks in which the target network selects the next best action while the current network estimates the $Q$-values given the action selected by the target. However, in most applications, it is unnecessary to estimate the value of each action choice. This is especially important in discrete problems with a large sets of possible actions where only a small portion of actions are suitable. For instance, in text summarization the output of the model is a vector of the distribution over the vocabulary and therefore, the output has the same dimension as the vocabulary size which

is usually selected to be between 50K to 150K. In most of the applications that use DDQN, the action space is limited to less than a few hundred. For instance, in an Atari game, the possible actions could be to move left, right, up, down, and shoot. Therefore, using DDQN would be easy for these types of applications. Recently, Wang *et al.* [152] proposed the idea of using a dueling net to overcome this problem. In their proposed method, rather than estimating the $Q$-values directly from the $Q$-net, two different values are estimated for the value function and advantage function as follows:

$$Q_\Psi(s_t, y_t) = V_\Psi(s_t) + A_\Psi(s_t, y_t) \qquad (31)$$

In order to be able to calculate the $V_\Psi(s_t)$, the value estimates is replicated $|\mathcal{A}|$ times. However, as discussed in [152], using Eq. (33) to calculate the $Q$ is bad and can potentially yield poor performance since Eq. (33) is unidentifiable in the sense that a constant can be added to $V_\Psi(s_t)$ and subtracted the same constant from $A_\Psi(s_t, y_t)$. To solve this problem, the authors suggested to force the advantage estimator to have a zero at the selected action:

$$Q_\Psi(s_t, y_t) = V_\Psi(s_t) + \Big(A_\Psi(s_t, y_t) - \max_y A_\Psi(s_t, y)\Big) \quad (32)$$

This way, for the action $y^* = \arg\max_y Q_\Psi(s_t, y) = \arg\max_y A_\Psi(s_t, y)$, $Q_\Psi(s_t, y^*) = V_\Psi(s_t)$ is obtained. As an alternative to Eq. (32) and to make the model more stable, the author suggested to replace the max operator with average:

$$Q_\Psi(s_t, y_t) = V_\Psi(s_t) + \Big(A_\Psi(s_t, y_t) - \frac{1}{|\mathcal{A}|}\sum_y A_\Psi(s_t, y)\Big)$$
$$(33)$$

Note that the dueling net will not decrease the number of actions but will provide a better normalization over the target distribution. Similar to DQN and DDQN, this model also suffers from the fact that there is no relation between the true values of $Q$-function and the estimation provided by the network. In Section V, we propose a simple and effective solution to overcome this problem by doing schedule sampling between the $Q$-value estimations and true $Q$-values to pre-train our function approximator. Fig. 4 summarizes some of the strengths and weaknesses of these different RL methods.

## IV. COMBINING RL WITH SEQ2SEQ MODELS

In this section, we will provide some of the recent models that combined the seq2seq training with reinforcement learning techniques. For most of these models, the main goal is to solve the train/test evaluation mismatch problem which exists in all previously described seq2seq models. This is usually done by adding a reward function to the training objective. There are a growing number of research works that used the REINFORCE algorithm to improve the current state-of-the-art seq2seq models. However, more advanced techniques such as Actor-Critic models, DQN and DDQN, have not been used often for these tasks. As mentioned earlier, one the main difficulties of using $Q$-Learning and its derivatives is the large action space for seq2seq models. For instance, in text summarization, the model should provide estimates for each word in the vocabulary and therefore the estimation

could be inferior even with a well trained model. Due to these reasons, researchers mostly focused on the easier yet problematic approaches such as REINFORCE algorithm to train the seq2seq model. Therefore, combining the power of $Q$-Learning training with seq2seq model is still considered to be an open area of research. Table III shows the policy, action, and reward function for each seq2seq task and Table IV summarizes these models along with the respective seq2seq application and specific RL algorithm they used to improve that application.

### A. Policy Gradient and REINFORCE Algorithm

As mentioned in Section III-A, in Policy Gradient (PG), the reward of the sampled sequence is observed at the end of the sequence generation and back-propagate that error equally to all the decoding steps according to Eq. (15). Also, we talked about the exposure bias problem that exists in seq2seq models during training the decoder because of using the Cross-Entropy (CE) error. The idea of improving generation by letting the model use its own predictions at training time was first proposed by Daume III *et al.* [153]. Based on their proposed method, SEARN, the structured prediction problems can be cast as a particular instance of reinforcement learning. The basic idea is to let the model use its own predictions at training time to produce a sequence of actions (e.g., the choice of the next word). Then, a greedy search algorithm is run to determine the optimal action at each time step, and the policy is trained to predict that action. An imitation learning framework was proposed by Ross *et al.* [154] in a method called DAGGER, where an oracle of the target word given the current predicted word is required. However, for tasks such as text summarization, computing the oracle is infeasible due to the large action space. This problem was later addressed by the '*Data As Demonstrator (DAD)*' model [155] where the target action at step $k$ is the $k^{th}$ action taken by the optimal policy. One drawback of DAD is that at every time step the target label is always selected from the ground-truth data and if the generated summaries are shorter than the ground-truth summaries, the model still forces to generate outputs that could already exist in the model. One way to avoid this problem in DAD is to use a method called End2EndBackProp [28] in which, at each step $t$, the top-$k$ actions are retrieved from the model and the normalized probabilities of these actions are used as weights (of importance) and the normalized combination of their representation is fed to the next decoding step.

Finally, REINFORCE algorithm [41] tries to overcome all these problems by using the PG rewarding function and avoiding the CE loss by using the sampled sequence as the ground-truth to train the seq2seq model, Eq. (18). In real-world applications, the training is usually started with the CE loss and a pre-trained model is acquired. Then, the REINFORCE algorithm is used to train the model. Some of the earliest adoptions of REINFORCE algorithm for training seq2seq models are in computer vision [89], [156], image captioning [19], and speech recognition [85]. Recently, other researchers showed that using a combination of CE loss and
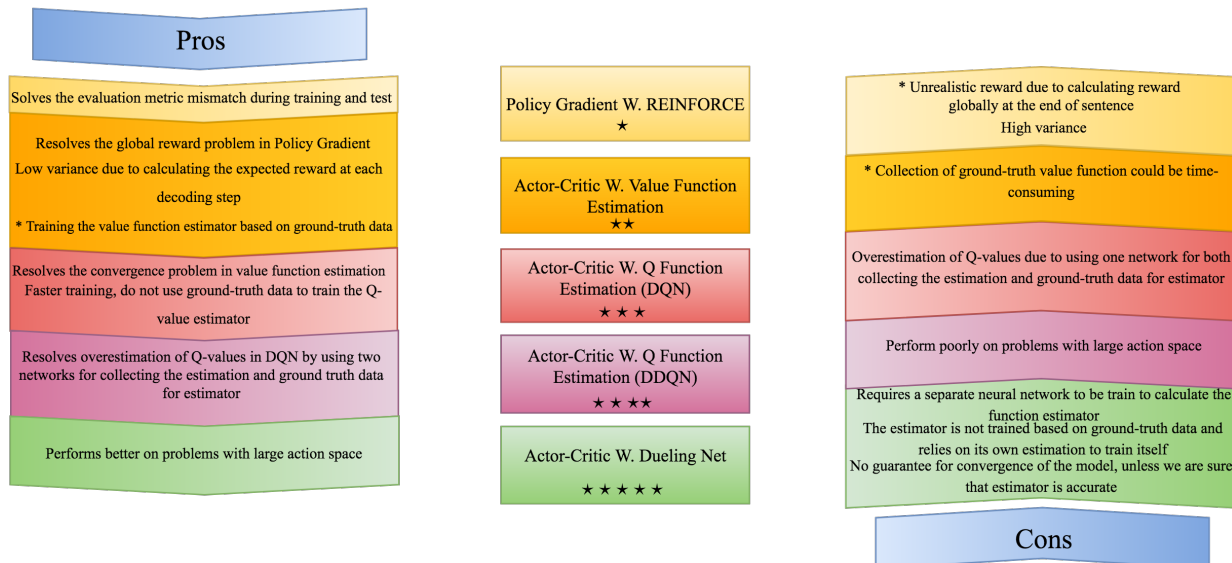
Fig. 4: A list of advantages and drawbacks of different RL models. The advantages are listed such that each method covers all the strengths of its previous methods and drawbacks are listed such that each method have all the weaknesses of the previous ones. For instance, Actor-Critic w. Dueling Net have all the pros of the previous models listed above it and Actor-Critic w. Value Function Estimation suffers from all the cons of the methods listed below it. The features that are also model-dependent are shown with '*' and those features do not exist in any other model. Each '⋆' shows how hard it is to implement these models in a real-world application.

TABLE III: Policy, Action, and Reward function for different seq2seq tasks.

| Seq2seq Task | Policy | Action | Reward |
|---|---|---|---|
| Text Summarization Headline Generation Machine Translation Question Generation | Attention-based models, pointer-generators, etc. | Selecting the next token for summary, headline, and translation | ROUGE, BLEU |
| Question Answering | Seq2seq model | Selecting the answer from a vocabulary or selecting the start and end index of the answer in the input document | F1 Score |
| Image Captioning Video Captioning | seq2seq model | Selecting the next token for the caption | CIDEr, SPICE, METEOR |
| Speech Recognition | Seq2seq model | Selecting the next token for the speech | Connectionist Temporal Classification (CTC) |
| Dialog Generation | Seq2seq model | Dialogue utterance to generate | BLEU Length of dialogue Diversity of dialogue |

REINFORCE loss could yield a better result than just simply performing the pre-training. In these models, the training is started by using the CE loss and is slowly switched from CE loss to REINFORCE loss to train the model. There are various ways in which one can do the transition from CE loss to REINFORCE loss. Ranzato *et al.* [28] used an incremental scheduling algorithm called 'MIXER' which combines DAG-GER [154] with DAD [155] methods. In this method, the RNN is trained with the cross-entropy loss for $N_{CE}$ epochs using the ground-truth sequences. This ensures that the model starts off with a much better policy than random because now the model can focus on promising regions of the search space. Then, they use an annealing schedule in order to gradually teach the model to produce stable sequences. Therefore, after the initial $N_{CE}$ epochs, they continue training the model for

$N_{CE} + N_R$ epochs, such that, for every sequence, they use the $\mathcal{L}_{CE}$ for the first $(T-\delta)$ steps, and the REINFORCE algorithm for the remaining $\delta$ steps. The MIXER model was successfully used on a variety of tasks such as text summarization, image captioning, and machine translation.

Another way to handle the transition from using CE loss to REINFORCE loss is to use the following combined loss:

$$\mathcal{L}_{mixed} = \eta\mathcal{L}_{REINFORCE} + (1-\eta)\mathcal{L}_{CE} \qquad (34)$$

where $\eta \in (0,1)$ is the parameter that controls the transition from CE to REINFORCE loss. In the beginning of the training, $\eta = 0$ and the model completely relies on CE loss, while as the training progresses, the $\eta$ value is increased in order to slowly reduce the effect of CE loss. By the end of the training process (where $\eta = 1$), the model completely uses the REINFORCE

loss for training. This mixed training loss was used in many of the recent works on text summarization [13], [46], [157], [158], paraphrase generation [159], image captioning [40], video captioning [160], speech recognition [161], dialogue generation [162], question answering [163], and question generation [62].

## B. Actor-Critic Models

One of the problems with the PG model is that we need to sample the full sequences of actions and observe the reward at the end of the generation. This, in general, will be problematic since the error of generation accumulates over time and usually for long sequences of actions, the final sequence is so far away from the ground-truth sequence. Thus, the reward of the final sequence would be small and model would take a lot of time to converge. To avoid this problem, Actor-Critic models observe the reward at each decoding step using the Critic model and fix the sequence of future actions that the Actor is going to take. The Critic model usually tries to maximize the advantage function through the estimation of value function or $Q$-function. As one of the early attempts of using AC models, Bahdanau *et al.* [39] and He *et al.* [164] used this model for the problem of machine translation. In [39], the author used temporal-difference (TD) learning for advantage function estimation by considering the $Q$-value for the next action, i.e., $Q(s_t, y_{t+1})$, as a surrogate for the its true value at time $t$, i.e., $V_\Psi(s_t)$. We mentioned that for a deterministic policy, $y^* = \arg\max_y Q(s, y)$, it follows that $Q(s, y^*) = V(s)$. Therefore, the $Q$-value for the next action could be used as the true estimates of the value function at current time. To accommodate for the large action space, they also use the shrinking estimation trick that was used in dueling net to push the estimate to be closer to their means. Additionally, the Critic training is done through the following mixed objective function:

$$\mathcal{L}(\Psi) = \frac{1}{2}\sum_i ||Q_\Psi(s_i, y_i) - q_i||^2 + \eta \bar{Q}_i$$
$$\bar{Q}_i = \sum_y \left( Q_\Psi(y, s_i) - \frac{1}{|\mathcal{A}|}\sum_{y'} Q_\Psi(y', s_i) \right) \quad (35)$$

where $q_i$ is the true estimation of $Q$ from a delayed Actor. The idea of using delayed Actor is similar to the idea used in Double $Q$-Learning where a delayed target network is used to get the estimation of the best action. Later, Zhang *et al.* [165] used a similar model on image captioning task.

He *et al.* [164] proposed a value network that uses a semantic matching and a context-coverage module and passed them through a dense layer to estimate the value function. However, their model requires a fully-trained seq2seq model to train the value network. Once the value network is trained, they use the trained seq2seq model and trained value estimation model to do the beam search during translation. Therefore, the value network is not used during the training of the seq2seq model. During inference, however, similar to the AlphaGo model [166], rather multiplying the advantage estimates (value or $Q$ estimates) to the policy probabilities (like in Eq. (23)), they combine the output of the seq2seq model and the value network as follows:

$$\eta \times \frac{1}{T}\log\pi(\hat{y}_{1\dots T}|X) + (1 - \eta) \times \log V_\Psi(\hat{y}_{1\dots T}) \quad (36)$$

where $V_\Psi(\hat{y}_{1\dots T})$ is the output of the value network and $\eta$ controls the effect of each score.

In a different model, Li *et al.* [167] proposed a model that controls the length of seq2seq model using RL-based ideas. They train a $Q$-value function approximator which estimates the future outcome of taking an action $y_t$ in the present and then incorporate it into a score $S(y_t)$ at each decoding step as follows:

$$S(y_t) = \log\pi(y_t|y_{t-1}, s_t) + \eta Q(X, y_{1\dots t}) \quad (37)$$

Specifically, the $Q$ function, in this work, takes only the hidden state at time $t$ and estimates the length of the remaining sequence. While decoding, they suggest an inference method that controls the length of the generated sequence as follows:

$$\hat{y}_t = \arg\max_y \log\pi(y|\hat{y}_{1\dots t-1}, X) - \eta||(T - t) - Q_\Psi(s_t)||^2 \quad (38)$$

Recently, Li *et al.* [46] proposed an AC model which uses a binary classifier as the Critic. In this specific model, the Critic tries to distinguish between the generated summary and the human-written summary via a neural network binary classifier. Once they pre-trained the Actor using CE loss, they start training the AC model alternatively using PG and the classifier score is considered as a surrogate for the value function. AC and PG were used also in the work of Liu *et al.* [135] where they combined AC and PG learning along with importance sampling to train a seq2seq model for image captioning. In this method, they used two different neural networks for $Q$-function estimation, i.e., $Q_\Psi$, and value estimation, i.e., $V_{\Psi'}$. They also used a mixed reward function that combines a weighted sum of $ROUGE$, $BLEU$, $METEOR$, and $CIDEr$ measures to achieve a higher performance on this task.

## C. Current RL-Based Model Issues

Throughout this paper, we discussed about various situations where using RL provides a better solution than traditional methods. However, utilizing RL methods creates its own training challenges and in most of the cases the improvement received from these models are not significant. In this section, we will discuss some of the issues that exist in current RL techniques used for seq2seq problems. As discussed in Section III, sample efficiency and high variance in RL models is one the of the main issues in applying them to seq2seq problems. Therefore, models such as RAML [173] and SPG [174] are proposed to provide a middle ground between the MLE and RL training. In RAML [173], a reward-aware perturbation is added to MLE while in SPG [174], the reward distribution is utilized for effective sampling of policy gradient. Recently, Tan *et al.* [175] provided a general formulation that connects the MLE and RL training through Entropy Regularized Policy Optimization (ERPO). However, even these solutions suffer from their own problems. RAML arguably suffers from the exposure bias while SPG requires a lot of engineering to work on a specific problem and, as shown in Table III, that is why REINFORCE-based models such as MIXER [28] are preferred in most of the current seq2seq problems.

TABLE IV: A summary of seq2seq applications that used various RL methods.

| Reference | Suffers From Exposure Bias | Mismatch on Train/Test Measure | Observe Full Reward | RL Algorithm | Seq2seq Application |
|---|---|---|---|---|---|
| Policy Gradient Based Models | | | | | |
| SEARN [153] | No | Yes | No Reward | PG | Sequence Labeling Syntactic Chunking |
| DAD [155] | No | Yes | No Reward | PG | Time-Series Modeling |
| Qin et al. [104] | No | Yes | Yes | PG | Relation Extraction |
| Yin et al. [108] | No | Yes | Yes | PG | Pronoun Resolution |
| MIXER [28] | No | No | Yes | PG w. REINFORCE | Machine Translation Text Summarization Image Captioning |
| Wu et al. [168] | No | No | Yes | PG w. REINFORCE | Text Summarization |
| Narayan et al. [158] | No | No | Yes | PG w. REINFORCE | Text Summarization |
| Kreutzer et al. [169] | No | No | Yes | PG w. REINFORCE | Machine Translation w. Human Bandit Feedback |
| Pan et al. [97] | No | No | Yes | PG w. REINFORCE | Natural Language Inference |
| Liang et al. [170] | No | No | Yes | PG w. REINFORCE | Semantic Parsing |
| Li et al. [162] | No | No | Yes | PG w. REINFORCE | Dialogue Generation |
| Yuan et al. [62] | No | No | Yes | PG w. REINFORCE | Question Generation |
| Mnih et al. [156] | Yes | No | Yes | PG w. REINFORCE | Computer Vision |
| Ba et al. [89] | Yes | No | Yes | PG w. REINFORCE | Computer Vision |
| Xu et al. [19] | Yes | No | Yes | PG w. REINFORCE | Image Captioning |
| Self-Critic Models with REINFORCE Algorithm | | | | | |
| Rennie et al. [40] | Yes | No | Yes | SC w. REINFORCE | Image Captioning |
| Paulus et al. [13] | No | No | Yes | SC w. REINFORCE | Text Summarization |
| Wang et al. [157] | No | No | Yes | SC w. REINFORCE | Text Summarization |
| Pasunuru et al. [160] | No | No | Yes | SC w. REINFORCE | Video Captioning |
| Yeung et al. [171] | No | No | Yes | SC w. REINFORCE | Action Detection in Video |
| Zhou et al. [161] | No | No | Yes | SC w. REINFORCE | Speech Recognition |
| Hu et al. [163] | No | No | Yes | SC w. REINFORCE | Question Answering |
| Actor-Critic Models with Policy Gradient and Q-Learning | | | | | |
| He et al. [164] | Yes | No | No | AC | Machine Translation |
| Li et al. [167] | Yes | No | No | AC | Machine Translation Text Summarization |
| Bahdanau et al. [39] | Yes | No | No | PG w. AC | Machine Translation |
| Li et al. [46] | Yes | No | No | PG w. AC | Text Summarization |
| Chen et al. [55] | Yes | No | Yes | PG w. AC | Text Summarization |
| Zhang et al. [165] | Yes | No | No | PG w. AC | Image Captioning |
| Liu et al. [135] | Yes | No | No | PG w. AC | Image Captioning |
| DARLA [172] | Yes | No | No | AC | Domain Adaptation |

Although REINFORCE-based models are simple to implement and provide better results, training these models is time-consuming and the improvement over baselines is usually marginal. This is why in most of the current works, these models are only used for fine-tuning purposes.

Aside from these issues, there are problems inherent to specific applications that make it hard for researchers to combine RL techniques with current seq2seq models. For instance, in most of the NLP problems, the output or action space is massive comparing to the size of actions in a robotic or game-playing problems. This is mostly due to the fact that in applications such as machine translation, text summarization, and image captioning the size of the output is equal the size of the vocabulary used during training. Now, compare this to an agent that plays a simple Atari game which requires deciding on usually less than 20 actions [176]. This will show the severity of this problem and the reward sparsity issue that exist in these applications.

Moreover, most of the current seq2seq models which use RL training, rely on well-defined reward function such as BLEU or ROUGE for providing feedback for the model. Although these are the standard metrics for evaluating various seq2seq models, relying on them creates a different set of problems. For instance, in abstractive text summarization, ROUGE and BLEU scores are being used as the standard metric for evaluation of summarization models. However, a good abstractive summary will definitely have a low ROUGE and BLEU score. This problem could be further investigated and possibly improved by Inverse Reinforcement Learning (IRL) [177] by forcing the model to learn its own rewarding function. However, to the best of our knowledge, no work has been done in this area.

Recently, new methods are introduced for game playing using RL algorithm which combine the best performing models in this area and apply some of the best practices used in previous models to achieve state-of-the-art results. Rainbow [178] and Quantile Nets [179] are among such frameworks. In Rainbow [178], the authors combine DDQN, prioritized experience buffer, dueling net, multi-step learning (using step-based reward rather than general reward), and distributional RL to achieve state-of-the-art in 57 games in the Atari 2600 framework. A similar ensembling method could also be useful to be applied for seq2seq tasks but this is also left for future works.

## V. RLSEQ2SEQ: AN OPEN-SOURCE LIBRARY FOR TRAINING SEQ2SEQ MODELS WITH RL METHODS

As part of this comprehensive study, we developed an open-source library which implements various RL techniques for the problem of abstractive text summarization. This library is made available at *www.github.com/yaserkl/RLSeq2Seq/*. Since experimenting each specific configuration of these models, even requires few days of training on GPUs, we encourage researchers, who use this library to build and enhance their own models, to also share their trained model at this website. In this section, we explain some of the important features of our library. As mentioned before, this library provides modules for abstractive text summarization. The core of our library is based on a model called pointer-generator [39] [12] which itself is based on Google TextSum model [40]. We also provide a similar imitation learning used in training REINFORCE algorithm to train the function approximator. This way, we propose training our DQN (DDQN, Dueling Net) using a schedule sampling in which we start training the model in the beginning based on ground-truth $Q$-values while as we move on with the training process, we completely rely on the function estimator to train the network. This could be seen as a pre-training step for the function approximator. Therefore, the model is guaranteed to start by using better ground-truth data since it is exposed to the true ground-truth values versus the random estimation it receives from the model itself. In summary, our library implements the following features:

- Adding temporal attention and intra-decoder attention that was proposed in [13].
- Adding scheduled sampling along with its differentiable relaxation proposed in [31] and E2EBackProb [28] for solving *exposure bias* problem.
- Adding adaptive training of REINFORCE algorithm by minimizing the mixed objective loss in Eq. (34).
- Providing Self-Critic training by adding the greedy reward as the baseline.
- Providing Actor-Critic training options for training the model using asynchronous training of Value Network, DQN, DDQN, and Dueling Net.
- Providing options for scheduled sampling for training of the $Q$-Function in DQN, DDQN, and Dueling Net.

### A. Experiments

To test the power of some of the studied models in this paper, we have done a range of various experiments using our open-source library. As mentioned in Section IV, most of the RL-based models play as a fine-tuning technique in seq2seq applications. Thus, we first pre-train our model for 15 epochs using only cross-entropy loss and then add the RL training for another 10 epochs. Our experiments follows the same setup as to the pointer-generator model [12] and we only show the results after activating the coverage mechanism. We activate the coverage mechanism only for the last epoch and select the best model using the evaluation data. We use

---

[39]https://github.com/abisee/pointer-generator

[40]https://github.com/tensorflow/models/tree/master/research/textsum

---

a linear scheduling probability as $\epsilon = step/totalsteps$ and we use $\epsilon = 1$ after activating the coverage so that the model completely relies on its own output for the rest of training and for E2EBackPropagation model, $K$ is set to 4. All experiments are done using two NVIDIA P100 GPUs one used for training the model and the other for select the best trained model based on the evaluation data.

*1) Analysis of the results:* Table V shows the results of our experiments based on ROUGE score on this dataset. All our ROUGE scores have a 95% confidence interval of at most $\pm 0.25$ as reported by the official ROUGE script. In this table, PG stands for Pointer-Generation and SS stands for Scheduled Sampling. As shown in this table, both scheduled sampling model and E2E model are superior to the pointer-generator. We have also used our framework to train the Self-Critic Policy Gradient (SCPG) based model proposed by Paulus *et al.* [13]. However, as shown in this table, although the SCPG improves the performance of the pointer-generator model, this improvement is very marginal. This result is totally in contrast with the result in the original paper [13] and shows that SCPG, as claimed by the authors, will not greatly improve the performance of the pointer-generator model. One of the main reason for this difference in the result of our experiment with the Paulus *et al.* [13] paper is that they use a completely different set of hyperparameters for training their model. For instance, the input for their encoder is 800 words while in our default setting, for all our experiments, it is set to 400. Also, the vocabulary size is set to 150K and 50K for input and output while our default is set to 50K for both input and output. Moreover, the size of hidden layers for encoder and decoder in their work is larger that our default values and they also use a pre-trained GloVe [180] word-embedding for training their model. Finally, we are comparing all these policy-gradient based models with an Actor-Critic model proposed by Chen *et al.* [55] which holds the state-of-the-art result in text summarization in CNN/DM dataset. As shown in Table V, this model is superior to any of the policy-gradient based models according to the ROUGE scores.

*2) Analysis of the training time:* In general, the pointer-generator framework requires more than 3 days of training for an effective results, while this time will also be expanded after adding the self-critic policy gradient. On average, each batch of training during MLE training will take 2-3 seconds while once we add the SCPG loss this time will be increased to 5-6 seconds which means the whole training time will be double after activation of the RL loss. On the other hand, the whole training time for the Actor-Critic model before and after RL activation is only a few hours which shows that not only it has superiority in the ROUGE score results but also the training process converges much faster than the other models.

## VI. CONCLUSION

In this paper, we provided a general overview of a specific type of deep learning models called sequence-to-sequence (seq2seq) models and discussed some of the recent advances in combining training of these models with Reinforcement Learning (RL) techniques. Seq2seq models are common in a

TABLE V: Analysis of ROUGE F1-Score after coverage and approximate amount of training time for various RL-based techniques.

| Method | ROUGE | | | Training Time |
|---|---|---|---|---|
| | 1 | 2 | L | |
| PG | 38.21 | 16.46 | 34.79 | 3-4 days |
| PG w. E2E | 38.24 | 16.48 | 34.97 | 3-4 days |
| PG w. SS Argmax-Sampling | 38.29 | 16.51 | 35.02 | 3-4 days |
| PG w. SS Argmax-Greedy | 38.65 | 16.77 | 35.37 | 3-4 days |
| PG w. SCPG | 38.77 | 16.98 | 35.32 | 6-7 days |
| Actor-Critic (Q-Learning) [55] | 40.88 | 17.80 | 38.54 | 3-4 hours |

wide range of applications from machine translation to speech recognition. However, traditional models in this field usually suffer from various problems during model training, such as inconsistency between the training objective and testing objective and *exposure bias*. Recently, with advances in deep reinforcement learning, researchers offered various types of solutions to combine the RL training with seq2seq training for alleviating the problems and challenges of training seq2seq models. In this paper, we summarized some of the most important works that tried to combine these two different techniques and provided an open-source library for the problem of abstractive text summarization that shows how one could train a seq2seq model using different RL techniques.

## REFERENCES

[1] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *NeurIPS*, 2014, pp. 3104–3112.

[2] S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer, "Scheduled sampling for sequence prediction with recurrent neural networks," in *NeurIPS*, 2015, pp. 1171–1179.

[3] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *EMNLP*, 2015, pp. 1412–1421.

[4] Y. Wu, M. Schuster, Z. Chen, Q. V. Le, M. Norouzi, W. Macherey, M. Krikun, Y. Cao, Q. Gao, K. Macherey *et al.*, "Google's neural machine translation system: Bridging the gap between human and machine translation," *arXiv preprint arXiv:1609.08144*, 2016.

[5] A. Vaswani, S. Bengio, E. Brevdo, F. Chollet, A. N. Gomez, S. Gouws, L. Jones, Ł. Kaiser, N. Kalchbrenner, N. Parmar *et al.*, "Tensor2tensor for neural machine translation," *arXiv preprint arXiv:1803.07416*, 2018.

[6] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv preprint arXiv:1409.0473*, 2014.

[7] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using rnn encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.

[8] S. Shen, Y. Cheng, Z. He, W. He, H. Wu, M. Sun, and Y. Liu, "Minimum risk training for neural machine translation," in *ACL*, vol. 1, 2016, pp. 1683–1692.

[9] A. M. Rush, S. Chopra, and J. Weston, "A neural attention model for abstractive sentence summarization," in *EMNLP*, 2015, pp. (379–389.

[10] S. Chopra, M. Auli, A. M. Rush, and S. Harvard, "Abstractive sentence summarization with attentive recurrent neural networks," in *NAACL-HLT*, 2016, pp. 93–98.

[11] R. Nallapati, B. Zhou, C. dos Santos, C. Gulcehre, and B. Xiang, "Abstractive text summarization using sequence-to-sequence rnns and beyond," in *SIGNLL*, 2016, pp. 280–290.

[12] A. See, P. J. Liu, and C. D. Manning, "Get to the point: Summarization with pointer-generator networks," in *ACL*, vol. 1, 2017, pp. 1073–1083.

[13] R. Paulus, C. Xiong, and R. Socher, "A deep reinforced model for abstractive summarization," *arXiv preprint arXiv:1705.04304*, 2017.

[14] R. Nallapati, F. Zhai, and B. Zhou, "Summarunner: A recurrent neural network based sequence model for extractive summarization of documents." in *AAAI*, 2017, pp. 3075–3081.

[15] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *ICASSP*. IEEE, 2013, pp. 6645–6649.

[16] D. Bahdanau, J. Chorowski, D. Serdyuk, P. Brakel, and Y. Bengio, "End-to-end attention-based large vocabulary speech recognition," in *ICASSP*. IEEE, 2016, pp. 4945–4949.

[17] D. Amodei, S. Ananthanarayanan, R. Anubhai, J. Bai, E. Battenberg, C. Case, J. Casper, B. Catanzaro, Q. Cheng, G. Chen *et al.*, "Deep speech 2: End-to-end speech recognition in english and mandarin," in *ICML*, 2016, pp. 173–182.

[18] S. Ö. Arık, M. Chrzanowski, A. Coates, G. Diamos, A. Gibiansky, Y. Kang, X. Li, J. Miller, A. Ng, J. Raiman *et al.*, "Deep voice: Real-time neural text-to-speech," in *ICML*, 2017, pp. 195–204.

[19] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, "Show, attend and tell: Neural image caption generation with visual attention," in *ICML*, 2015, pp. 2048–2057.

[20] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and tell: A neural image caption generator," in *CVPR*. IEEE, 2015, pp. 3156–3164.

[21] A. Karpathy and L. Fei-Fei, "Deep visual-semantic alignments for generating image descriptions," in *CVPR*, 2015, pp. 3128–3137.

[22] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[23] J. L. Elman, "Finding structure in time," *Cognitive science*, vol. 14, no. 2, pp. 179–211, 1990.

[24] C.-Y. LIN, "Rouge: A package for automatic evaluation of summaries," in *Proc. of Workshop on Text Summarization Branches Out*, 2004.

[25] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *ACL*, 2002, pp. 311–318.

[26] S. Banerjee and A. Lavie, "Meteor: An automatic metric for mt evaluation with improved correlation with human judgments," in *Proceedings of the ACL workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005, pp. 65–72.

[27] R. Vedantam, C. Lawrence Zitnick, and D. Parikh, "Cider: Consensus-based image description evaluation," in *CVPR*, 2015, pp. 4566–4575.

[28] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba, "Sequence level training with recurrent neural networks," *arXiv preprint arXiv:1511.06732*, 2015.

[29] J. Su, J. Xu, X. Qiu, and X. Huang, "Incorporating discriminator in sentence generation: a gibbs sampling method," *arXiv preprint arXiv:1802.08970*, 2018.

[30] F. Huszár, "How (not) to train your generative model: Scheduled sampling, likelihood, adversary?" *arXiv preprint arXiv:1511.05101*, 2015.

[31] K. Goyal, C. Dyer, and T. Berg-Kirkpatrick, "Differentiable scheduled sampling for credit assignment," in *ACL*, vol. 2, 2017, pp. 366–371.

[32] L. Yu, W. Zhang, J. Wang, and Y. Yu, "Seqgan: sequence generative adversarial nets with policy gradient," in *AAAI*, vol. 31, 2017, pp. 2852–2858.

[33] K. Lin, D. Li, X. He, Z. Zhang, and M.-T. Sun, "Adversarial ranking for language generation," in *NeurIPS*, 2017, pp. 3155–3165.

[34] J. Guo, S. Lu, H. Cai, W. Zhang, Y. Yu, and J. Wang, "Long text generation via adversarial training with leaked information," *arXiv preprint arXiv:1709.08624*, 2017.

[35] T. Che, Y. Li, R. Zhang, R. D. Hjelm, W. Li, Y. Song, and Y. Bengio, "Maximum-likelihood augmented discrete generative adversarial networks," *arXiv preprint arXiv:1702.07983*, 2017.

[36] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *NeurIPS*, 2014, pp. 2672–2680.

[37] Y. Zhang, Z. Gan, K. Fan, Z. Chen, R. Henao, D. Shen, and L. Carin, "Adversarial feature matching for text generation," in *ICML*, 2017, pp. 4006–4015.

[38] Z. Shi, X. Chen, X. Qiu, and X. Huang, "Towards diverse text generation with inverse reinforcement learning," *arXiv preprint arXiv:1804.11258*, 2018.
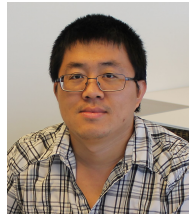
[39] D. Bahdanau, P. Brakel, K. Xu, A. Goyal, R. Lowe, J. Pineau, A. Courville, and Y. Bengio, "An actor-critic algorithm for sequence prediction," in *ICLR*, 2017.

[40] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," in *CVPR*, 2017, pp. 7008–7024.

[41] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," in *Reinforcement Learning*. Springer, 1992, pp. 5–32.

[42] R. S. Sutton and A. G. Barto, "Reinforcement learning: An introduction," *Cambridge: MIT press*, vol. 1, no. 1, 1998.

[43] S. Levine, C. Finn, T. Darrell, and P. Abbeel, "End-to-end training of deep visuomotor policies," *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334–1373, 2016.

[44] I. Lenz, H. Lee, and A. Saxena, "Deep learning for detecting robotic grasps," *The International Journal of Robotics Research*, vol. 34, no. 4-5, pp. 705–724, 2015.

[45] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski *et al.*, "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, p. 529, 2015.

[46] P. Li, L. Bing, and W. Lam, "Actor-critic based training framework for abstractive summarization," *arXiv preprint arXiv:1803.11070*, 2018.

[47] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath, "A brief survey of deep reinforcement learning," *arXiv preprint arXiv:1708.05866*, 2017.

[48] Y. Li, "Deep reinforcement learning: An overview," *arXiv preprint arXiv:1701.07274*, 2017.

[49] S. Team, "Deep learning for siri's voice: On-device deep mixture density networks for hybrid unit selection synthesis," *Apple Machine Learning Journal*, vol. 1, no. 4, 2017.

[50] G. Klein, Y. Kim, Y. Deng, J. Senellart, and A. Rush, "Opennmt: Open-source toolkit for neural machine translation," *Proceedings of ACL, System Demonstrations*, pp. 67–72, 2017.

[51] M. X. Chen, O. Firat, A. Bapna, M. Johnson, W. Macherey, G. Foster, L. Jones, N. Parmar, M. Schuster, Z. Chen *et al.*, "The best of both worlds: Combining recent advances in neural machine translation," in *ACL*, 2018, pp. 76–86.

[52] J. Ling and A. Rush, "Coarse-to-fine attention models for document summarization," in *Proceedings of the Workshop on New Frontiers in Summarization*, 2017, pp. 33–42.

[53] Q. Zhou, N. Yang, F. Wei, and M. Zhou, "Selective encoding for abstractive sentence summarization," in *ACL*, vol. 1, 2017, pp. 1095–1104.

[54] J. Tan, X. Wan, and J. Xiao, "Abstractive document summarization with a graph-based attentional neural model," in *ACL*, vol. 1, 2017, pp. 1171–1181.

[55] Y.-C. Chen and M. Bansal, "Fast abstractive summarization with reinforce-selected sentence rewriting," in *ACL*, vol. 1, 2018, pp. 675–686.

[56] J. Lin, X. Sun, S. Ma, and Q. Su, "Global encoding for abstractive summarization," in *ACL*, 2018, pp. 163–169.

[57] W.-T. Hsu, C.-K. Lin, M.-Y. Lee, K. Min, J. Tang, and M. Sun, "A unified model for extractive and abstractive summarization using inconsistency loss," in *ACL*, 2018, pp. 132–141.

[58] Y. Xia, F. Tian, L. Wu, J. Lin, T. Qin, N. Yu, and T.-Y. Liu, "Deliberation networks: Sequence generation beyond one-pass decoding," in *NeurIPS*, 2017, pp. 1782–1792.

[59] I. V. Serban, A. García-Durán, C. Gulcehre, S. Ahn, S. Chandar, A. Courville, and Y. Bengio, "Generating factoid questions with recurrent neural networks: The 30m factoid question-answer corpus," in *ACL*, vol. 1, 2016, pp. 588–598.

[60] N. Mostafazadeh, I. Misra, J. Devlin, M. Mitchell, X. He, and L. Vanderwende, "Generating natural questions about an image," in *ACL*, vol. 1, 2016, pp. 1802–1813.

[61] Z. Yang, J. Hu, R. Salakhutdinov, and W. Cohen, "Semi-supervised qa with generative domain-adaptive nets," in *ACL*, vol. 1, 2017, pp. 1040–1050.

[62] X. Yuan, T. Wang, C. Gulcehre, A. Sordoni, P. Bachman, S. Zhang, S. Subramanian, and A. Trischler, "Machine comprehension by text-to-text neural question generation," in *Proceedings of the 2nd Workshop on Representation Learning for NLP*, 2017, pp. 15–25.

[63] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh, "VQA: Visual question answering," in *ICCV*, 2015, pp. 2425–2433.

[64] C. Xiong, V. Zhong, and R. Socher, "Dynamic coattention networks for question answering," in *ICLR*, 2016.

[65] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, "Stacked attention networks for image question answering," in *CVPR*, 2016, pp. 21–29.

[66] O. Vinyals and Q. Le, "A neural conversational model," *arXiv preprint arXiv:1506.05869*, 2015.

[67] J. Li, M. Galley, C. Brockett, J. Gao, and B. Dolan, "A diversity-promoting objective function for neural conversation models," in *NAACL-HLT*, 2016, pp. 110–119.

[68] I. V. Serban, A. Sordoni, Y. Bengio, A. C. Courville, and J. Pineau, "Building end-to-end dialogue systems using generative hierarchical neural network models." in *AAAI*, vol. 16, 2016, pp. 3776–3784.

[69] A. Bordes, Y.-L. Boureau, and J. Weston, "Learning end-to-end goal-oriented dialog," *arXiv preprint arXiv:1605.07683*, 2016.

[70] V. Zhong, C. Xiong, and R. Socher, "Seq2SQL: Generating structured queries from natural language using reinforcement learning," 2018. [Online]. Available: https://openreview.net/forum?id=Syx6bz-Ab

[71] X. Xu, C. Liu, and D. Song, "Sqlnet: Generating structured queries from natural language without reinforcement learning," *arXiv preprint arXiv:1711.04436*, 2017.

[72] R. Kiros, R. Salakhutdinov, and R. Zemel, "Multimodal neural language models," in *ICML*, 2014, pp. 595–603.

[73] R. Kiros, R. Salakhutdinov, and R. S. Zemel, "Unifying visual-semantic embeddings with multimodal neural language models," *arXiv preprint arXiv:1411.2539*, 2014.

[74] X. Chen and C. L. Zitnick, "Mind's eye: A recurrent visual representation for image caption generation," in *CVPR*. IEEE, 2015, pp. 2422–2431.

[75] J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, and A. Yuille, "Deep captioning with multimodal recurrent neural networks (m-rnn)," *arXiv preprint arXiv:1412.6632*, 2014.

[76] H. Fang, S. Gupta, F. Iandola, R. K. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. C. Platt *et al.*, "From captions to visual concepts and back," in *CVPR*. IEEE, 2015, pp. 1473–1482.

[77] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell, "Long-term recurrent convolutional networks for visual recognition and description," in *CVPR*, 2015, pp. 2625–2634.

[78] S. Venugopalan, H. Xu, J. Donahue, M. Rohrbach, R. Mooney, and K. Saenko, "Translating videos to natural language using deep recurrent neural networks," in *NAACL-HLT*, 2015, pp. 1494–1504.

[79] S. Venugopalan, M. Rohrbach, J. Donahue, R. Mooney, T. Darrell, and K. Saenko, "Sequence to sequence-video to text," in *ICCV*, 2015, pp. 4534–4542.

[80] S. Venugopalan, L. A. Hendricks, R. Mooney, and K. Saenko, "Improving lstm-based video description with linguistic knowledge mined from text," in *EMNLP*, 2016, pp. 1961–1966.

[81] P. Sermanet, D. Eigen, X. Zhang, M. Mathieu, R. Fergus, and Y. Lecun, "Overfeat: Integrated recognition, localization and detection using convolutional networks," in *ICLR*, 2014.

[82] D. Erhan, C. Szegedy, A. Toshev, and D. Anguelov, "Scalable object detection using deep neural networks," in *CVPR*, 2014, pp. 2147–2154.

[83] R. Girshick, "Fast r-cnn," *arXiv preprint arXiv:1504.08083*, 2015.

[84] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *NeurIPS*, 2015, pp. 91–99.

[85] A. Graves and N. Jaitly, "Towards end-to-end speech recognition with recurrent neural networks," in *ICML*, 2014, pp. 1764–1772.

[86] Y. Miao, M. Gowayyed, and F. Metze, "Eesen: End-to-end speech recognition using deep rnn models and wfst-based decoding," in *IEEE Workshop on Automatic Speech Recognition and Understanding (ASRU)*. IEEE, 2015, pp. 167–174.

[87] H. Ze, A. Senior, and M. Schuster, "Statistical parametric speech synthesis using deep neural networks," in *ICASSP*. IEEE, 2013, pp. 7962–7966.

[88] Y. Fan, Y. Qian, F.-L. Xie, and F. K. Soong, "Tts synthesis with bidirectional lstm based recurrent neural networks," in *ISCA*, 2014.

[89] J. Ba, V. Mnih, and K. Kavukcuoglu, "Multiple object recognition with visual attention," *arXiv preprint arXiv:1412.7755*, 2014.

[90] O. Vinyals, M. Fortunato, and N. Jaitly, "Pointer networks," in *NeurIPS*, 2015, pp. 2692–2700.

[91] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017, pp. 6000–6010.

[92] A. Radford, R. Jozefowicz, and I. Sutskever, "Learning to generate reviews and discovering sentiment," *arXiv preprint arXiv:1704.01444*, 2017.

[93] C. dos Santos and M. Gatti, "Deep convolutional neural networks for sentiment analysis of short texts," in *COLING*, 2014, pp. 69–78.

[94] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, "Recursive deep models for semantic compositionality over a sentiment treebank," in *EMNLP*, 2013, pp. 1631–1642.

[95] A. Conneau, D. Kiela, H. Schwenk, L. Barrault, and A. Bordes, "Supervised learning of universal sentence representations from natural language inference data," in *EMNLP*, 2017, pp. 670–680.

[96] S. Kim, J.-H. Hong, I. Kang, and N. Kwak, "Semantic sentence matching with densely-connected recurrent and co-attentive information," *arXiv preprint arXiv:1805.11360*, 2018.

[97] B. Pan, Y. Yang, Z. Zhao, Y. Zhuang, D. Cai, and X. He, "Discourse marker augmented network with reinforcement learning for natural language inference," in *ACL*, vol. 1, 2018, pp. 989–999.

[98] N. FitzGerald, O. Täckström, K. Ganchev, and D. Das, "Semantic role labeling with neural network factors," in *EMNLP*, 2015, pp. 960–970.

[99] L. He, M. Lewis, and L. Zettlemoyer, "Question-answer driven semantic role labeling: Using natural language to annotate natural language," in *EMNLP*, 2015, pp. 643–653.

[100] D. Marcheggiani, A. Frolov, and I. Titov, "A simple and accurate syntax-agnostic neural model for dependency-based semantic role labeling," in *CoNLL*, 2017, pp. 411–420.

[101] D. Marcheggiani and I. Titov, "Encoding sentences with graph convolutional networks for semantic role labeling," in *EMNLP*, 2017, pp. 1506–1515.

[102] M. Mintz, S. Bills, R. Snow, and D. Jurafsky, "Distant supervision for relation extraction without labeled data," in *ACL*, 2009, pp. 1003–1011.

[103] X. Huang *et al.*, "Attention-based convolutional neural network for semantic relation extraction," in *ICLR*, 2016, pp. 2526–2536.

[104] P. Qin, W. Xu, and W. Y. Wang, "Robust distant supervision relation extraction via deep reinforcement learning," in *ACL*, vol. 1, 2018, pp. 2137–2147.

[105] C. Chen and V. Ng, "Chinese zero pronoun resolution with deep neural networks," in *ACL*, vol. 1, 2016, pp. 778–788.

[106] Q. Yin, Y. Zhang, W. Zhang, and T. Liu, "Chinese zero pronoun resolution with deep memory network," in *EMNLP*, 2017, pp. 1309–1318.

[107] T. H. Trinh and Q. V. Le, "A simple method for commonsense reasoning," *arXiv preprint arXiv:1806.02847*, 2018.

[108] Q. Yin, Y. Zhang, W. Zhang, T. Liu, and W. Y. Wang, "Deep reinforcement learning for chinese zero pronoun resolution," in *ACL*, 2018, pp. 569–578.

[109] P. Anderson, B. Fernando, M. Johnson, and S. Gould, "Spice: Semantic propositional image caption evaluation," in *ECCV*, 2016.

[110] A. Axelrod, X. He, and J. Gao, "Domain adaptation via pseudo in-domain data selection," in *EMNLP*, 2011, pp. 355–362.

[111] K. M. Hermann, T. Kocisky, E. Grefenstette, L. Espeholt, W. Kay, M. Suleyman, and P. Blunsom, "Teaching machines to read and comprehend," in *NeurIPS*, 2015, pp. 1693–1701.

[112] M. Grusky, M. Naaman, and Y. Artzi, "Newsroom: A dataset of 1.3 million summaries with diverse extractive strategies," in *NAACL-HLT*, New Orleans, Louisiana, June 2018. [Online]. Available: https://summari.es/newsroom.pdf

[113] D. Graff, J. Kong, K. Chen, and K. Maeda, "English gigaword," *Linguistic Data Consortium, Philadelphia*, vol. 4, p. 1, 2003.

[114] P. Rajpurkar, J. Zhang, K. Lopyrev, and P. Liang, "Squad: 100,000+ questions for machine comprehension of text," *arXiv preprint arXiv:1606.05250*, 2016.

[115] P. Rajpurkar, R. Jia, and P. Liang, "Know what you don't know: Unanswerable questions for squad," *arXiv preprint arXiv:1806.03822*, 2018.

[116] M. Joshi, E. Choi, D. Weld, and L. Zettlemoyer, "Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension," in *ACL*, vol. 1, 2017, pp. 1601–1611.

[117] J. Tiedemann, "News from opusa collection of multilingual parallel corpora with tools and interfaces 1 index of subjects and terms 13." [Online]. Available: http://opus.nlpl.eu/

[118] J. Dodge, A. Gane, X. Zhang, A. Bordes, S. Chopra, A. Miller, A. Szlam, and J. Weston, "Evaluating prerequisite qualities for learning end-to-end dialog systems," *arXiv preprint arXiv:1511.06931*, 2015.

[119] C. Danescu-Niculescu-Mizil and L. Lee, "Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs," in *Proceedings of the 2nd Workshop on Cognitive Modeling and Computational Linguistics*, 2011, pp. 76–87.

[120] P. Pasupat and P. Liang, "Compositional semantic parsing on semi-structured tables," in *ACL*, 2015, pp. 1470–1480.

[121] Y. Wang, J. Berant, and P. Liang, "Building a semantic parser overnight," in *ACL*, vol. 1, 2015, pp. 1332–1342.

[122] J. McAuley, R. Pandey, and J. Leskovec, "Inferring networks of substitutable and complementary products," in *KDD*. ACM, 2015, pp. 785–794.

[123] S. R. Bowman, G. Angeli, C. Potts, and C. D. Manning, "A large annotated corpus for learning natural language inference," in *EMNLP*, 2015.

[124] A. Williams, N. Nangia, and S. Bowman, "A broad-coverage challenge corpus for sentence understanding through inference," in *NAACL-HLT*, vol. 1, 2018, pp. 1112–1122.

[125] M. Palmer, D. Gildea, and P. Kingsbury, "The proposition bank: An annotated corpus of semantic roles," *Computational linguistics*, vol. 31, no. 1, pp. 71–106, 2005.

[126] J. R. Finkel, T. Grenager, and C. Manning, "Incorporating non-local information into information extraction systems by gibbs sampling," in *ACL*, 2005, pp. 363–370.

[127] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European conference on computer vision*. Springer, 2014, pp. 740–755.

[128] V. Ordonez, G. Kulkarni, and T. L. Berg, "Im2text: Describing images using 1 million captioned photographs," in *NeurIPS*, 2011, pp. 1143–1151.

[129] J. Xu, T. Mei, T. Yao, and Y. Rui, "Msr-vtt: A large video description dataset for bridging video and language," in *CVPR*. IEEE, 2016, pp. 5288–5296.

[130] D. L. Chen and W. B. Dolan, "Collecting highly parallel data for paraphrase evaluation," in *ACL*, 2011, pp. 190–200.

[131] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.

[132] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, "Librispeech: an asr corpus based on public domain audio books," in *ICASSP*. IEEE, 2015, pp. 5206–5210.

[133] W. Zaremba and I. Sutskever, "Reinforcement learning neural turing machines-revised," *arXiv preprint arXiv:1505.00521*, 2015.

[134] T. Jie and P. Abbeel, "On a connection between importance sampling and the likelihood ratio policy gradient," in *NeurIPS*, 2010, pp. 1000–1008.

[135] S. Liu, Z. Zhu, N. Ye, S. Guadarrama, and K. Murphy, "Improved image captioning via policy gradient optimization of spider," in *ICCV*, vol. 3, 2017.

[136] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *ICML*, 2016, pp. 1928–1937.

[137] R. Munos, T. Stepleton, A. Harutyunyan, and M. Bellemare, "Safe and efficient off-policy reinforcement learning," in *NeurIPS*, 2016, pp. 1054–1062.

[138] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu, "Reinforcement learning with unsupervised auxiliary tasks," *arXiv preprint arXiv:1611.05397*, 2016.

[139] A. Gruslys, M. G. Azar, M. G. Bellemare, and R. Munos, "The reactor: A sample-efficient actor-critic architecture," *arXiv preprint arXiv:1704.04651*, 2017.

[140] Z. Wang, V. Bapst, N. Heess, V. Mnih, R. Munos, K. Kavukcuoglu, and N. de Freitas, "Sample efficient actor-critic with experience replay," in *ICLR*, 2017.

[141] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *NeurIPS*, 2000, pp. 1057–1063.

[142] G. Tucker, S. Bhupatiraju, S. Gu, R. E. Turner, Z. Ghahramani, and S. Levine, "The mirage of action-dependent baselines in reinforcement learning," *ICLR Workshop*, 2018.

[143] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel, "High-dimensional continuous control using generalized advantage estimation," in *ICLR*, 2016.

[144] R. S. Sutton and A. G. Barto, "Reinforcement learning: An introduction," *MIT press*, 2018.

[145] C. J. Watkins and P. Dayan, "Q-learning," *Machine learning*, vol. 8, no. 3-4, pp. 279–292, 1992.

[146] J. N. Tsitsiklis, "Asynchronous stochastic approximation and q-learning," *Machine learning*, vol. 16, no. 3, pp. 185–202, 1994.

[147] L.-J. Lin, "Self-improving reactive agents based on reinforcement learning, planning and teaching," *Machine learning*, vol. 8, no. 3-4, pp. 293–321, 1992.

[148] T. Schaul, J. Quan, I. Antonoglou, and D. Silver, "Prioritized experience replay," *arXiv preprint arXiv:1511.05952*, 2015.

[149] H. V. Hasselt, "Double q-learning," in *NeurIPS*, 2010, pp. 2613–2621.

[150] S. Thrun and A. Schwartz, "Issues in using function approximation for reinforcement learning," in *Proceedings of the 1993 Connectionist Models Summer School Hillsdale, NJ. Lawrence Erlbaum*, 1993.

[151] H. van Hasselt, A. Guez, and D. Silver, "Deep reinforcement learning with double q-learning," in *AAAI*, 2016.

[152] Z. Wang, T. Schaul, M. Hessel, H. Hasselt, M. Lanctot, and N. Freitas, "Dueling network architectures for deep reinforcement learning," in *ICML*, 2016, pp. 1995–2003.

[153] H. Daumé, J. Langford, and D. Marcu, "Search-based structured prediction," *Machine learning*, vol. 75, no. 3, pp. 297–325, 2009.

[154] S. Ross, G. Gordon, and D. Bagnell, "A reduction of imitation learning and structured prediction to no-regret online learning," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, 2011, pp. 627–635.

[155] A. Venkatraman, M. Hebert, and J. A. Bagnell, "Improving multi-step prediction of learned time series models." in *AAAI*, 2015, pp. 3024–3030.

[156] V. Mnih, N. Heess, A. Graves *et al.*, "Recurrent models of visual attention," in *NeurIPS*, 2014, pp. 2204–2212.

[157] L. Wang, J. Yao, Y. Tao, L. Zhong, W. Liu, and Q. Du, "A reinforced topic-aware convolutional sequence-to-sequence model for abstractive text summarization," *arXiv preprint arXiv:1805.03616*, 2018.

[158] S. Narayan, S. B. Cohen, and M. Lapata, "Ranking sentences for extractive summarization with reinforcement learning," in *NAACL-HLT*, vol. 1, 2018, pp. 1747–1759.

[159] Z. Li, X. Jiang, L. Shang, and H. Li, "Paraphrase generation with deep reinforcement learning," *arXiv preprint arXiv:1711.00279*, 2017.

[160] R. Pasunuru and M. Bansal, "Reinforced video captioning with entailment rewards," in *EMNLP*, 2017, pp. 979–985.

[161] Y. Zhou, C. Xiong, and R. Socher, "Improving end-to-end speech recognition with policy learning," *arXiv preprint arXiv:1712.07101*, 2017.

[162] J. Li, W. Monroe, A. Ritter, D. Jurafsky, M. Galley, and J. Gao, "Deep reinforcement learning for dialogue generation," in *EMNLP*, 2016, pp. 1192–1202.

[163] M. Hu, Y. Peng, and X. Qiu, "Reinforced mnemonic reader for machine comprehension," *CoRR, abs/1705.02798*, 2017.

[164] D. He, H. Lu, Y. Xia, T. Qin, L. Wang, and T. Liu, "Decoding with value networks for neural machine translation," in *NeurIPS*, 2017, pp. 177–186.

[165] L. Zhang, F. Sung, F. Liu, T. Xiang, S. Gong, Y. Yang, and T. M. Hospedales, "Actor-critic sequence training for image captioning," *arXiv preprint arXiv:1706.09601*, 2017.

[166] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.

[167] J. Li, W. Monroe, and D. Jurafsky, "Learning to decode for future success," *arXiv preprint arXiv:1701.06549*, 2017.

[168] Y. Wu and B. Hu, "Learning to Extract Coherent Summary via Deep Reinforcement Learning," *arXiv preprint arXiv:1804.07036*, 2018.

[169] J. Kreutzer, J. Uyheng, and S. Riezler, "Reliability and learnability of human bandit feedback for sequence-to-sequence reinforcement learning," in *ACL*, vol. 1, 2018, pp. 1777–1788.

[170] C. Liang, J. Berant, Q. Le, K. D. Forbus, and N. Lao, "Neural symbolic machines: Learning semantic parsers on freebase with weak supervision," in *ACL*, vol. 1, 2017, pp. 23–33.

[171] S. Yeung, O. Russakovsky, G. Mori, and L. Fei-Fei, "End-to-end learning of action detection from frame glimpses in videos," in *CVPR*, 2016, pp. 2678–2687.

[172] I. Higgins, A. Pal, A. Rusu, L. Matthey, C. Burgess, A. Pritzel, M. Botvinick, C. Blundell, and A. Lerchner, "Darla: Improving zero-shot transfer in reinforcement learning," in *International Conference on Machine Learning*, 2017, pp. 1480–1490.

[173] M. Norouzi, S. Bengio, N. Jaitly, M. Schuster, Y. Wu, D. Schuurmans *et al.*, "Reward augmented maximum likelihood for neural structured prediction," in *NeurIPS*, 2016, pp. 1723–1731.

[174] N. Ding and R. Soricut, "Cold-start reinforcement learning with softmax policy gradient," in *NeurIPS*, 2017, pp. 2817–2826.

[175] B. Tan, Z. Hu, Z. Yang, R. Salakhutdinov, and E. Xing, "Connecting the dots between mle and rl for sequence generation," *arXiv preprint arXiv:1811.09740*, 2018.

[176] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller, "Playing atari with deep reinforcement learning," *arXiv preprint arXiv:1312.5602*, 2013.

[177] B. D. Ziebart, A. Maas, J. A. Bagnell, and A. K. Dey, "Maximum entropy inverse reinforcement learning," in *AAAI*. AAAI Press, 2008, pp. 1433–1438.

[178] M. Hessel, J. Modayil, H. Van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver, "Rainbow: Combining improvements in deep reinforcement learning," *arXiv preprint arXiv:1710.02298*, 2017.

[179] W. Dabney, G. Ostrovski, D. Silver, and R. Munos, "Implicit quantile networks for distributional reinforcement learning," *arXiv preprint arXiv:1806.06923*, 2018.

[180] J. Pennington, R. Socher, and C. Manning, "Glove: Global vectors for word representation," in *EMNLP*, 2014, pp. 1532–1543.

**Yaser Keneshloo** received his Masters degree in Computer Engineering from Iran University of Science and Technology in 2012. Currently, he is a Ph.D candidate in the Department of Computer Science at Virginia Tech. His research interests includes machine learning, data mining, and deep learning.

**Tian Shi** Tian Shi received the Ph.D. degree in Physical Chemistry from Wayne State University in 2016. He is working toward the Ph.D. degree in the Department of Computer Science, Virginia Tech. His research interests include data mining, deep learning, topic modeling, and text summarization.

**Naren Ramakrishnan** is the Thomas L. Phillips Professor of Engineering at Virginia Tech. He directs the Discovery Analytics Center, a university-wide effort that brings together researchers from computer science, statistics, mathematics, and electrical and computer engineering to tackle knowledge discovery problems in important areas of national interest. His work has been featured in the Wall Street Journal, Newsweek, Smithsonian Magazine, PBS/NoVA Next, Chronicle of Higher Education, and Popular Science, among other venues. Ramakrishnan serves on the editorial boards of IEEE Computer, ACM Transactions on Knowledge Discovery from Data, Data Mining and Knowledge Discovery, IEEE Transactions on Knowledge and Data Engineering, and other journals. He received his PhD in Computer Sciences from Purdue University.

**Chandan K. Reddy** is an Associate Professor in the Department of Computer Science at Virginia Tech. He received his Ph.D. from Cornell University and M.S. from Michigan State University. His primary research interests are Data Mining and Machine Learning with applications to Healthcare Analytics and Social Network Analysis. His research is funded by the National Science Foundation, the National Institutes of Health, the Department of Transportation, and the Susan G. Komen for the Cure Foundation. He has published over 95 peer-reviewed articles in leading conferences and journals. He received several awards for his research work including the Best Application Paper Award at ACM SIGKDD conference in 2010, Best Poster Award at IEEE VAST conference in 2014, Best Student Paper Award at IEEE ICDM conference in 2016, and was a finalist of the INFORMS Franz Edelman Award Competition in 2011. He is an associate editor of the ACM Transactions on Knowledge Discovery and Data Mining and PC Co-Chair of ASONAM 2018. He is a senior member of the IEEE and life member of the ACM.